

# Applications and Industry®

UNIVERSITY OF HAWAII  
LIBRARY

JUN 3 8 38 AM '70

November 1961



## Transactions Papers

### General Applications Division

- 60-1028 Effect of Variable Plasma Conductivity—MHD Converter.....Coe, Eisen . . . 225
- 61-208 Economic Justification of Railway Electrification in U. S.....Cross . . . 232

### Industry Division

- 61-802 Optimum Synthesis of Random Sampling Multipole Filters.....Hsieh . . . 239
- 61-743 Sampled-Data Control Systems with Transport Lag.....Šiljak . . . 247
- 61-819 Mathematical Models for Time-Domain Design.....Wang . . . 252
- 61-712 Signal Stabilization of Control System.....Oldenburger, Sridhar . . . 260
- 61-710 Frequency Response of Nonlinear Closed-Loop Systems.....McAllister . . . 268
- 61-709 Stabilization of Feedback Systems.....Mahalanabis . . . 277
- 60-1066 Thermoelectricity Application Considerations.....Sorensen . . . 285
- 61-713 Advances in Analysis and Synthesis of Nonlinear Systems.....Wolf . . . 289
- 61-752 Feedback Compensation: A Design Technique..Thaler, Bronzino, Kirk . . . 300
- 61-711 Can Electric Actuators Meet Missile Requirements?....Newton, Rasche . . . 306
- Conference Papers Open for Discussion.....See 3rd Cover

© Copyright 1961 by American Institute of Electrical Engineers

NUMBER 57

*Published Bimonthly by*

AMERICAN INSTITUTE OF ELECTRICAL ENGINEERS

Instrumentation Division

61-721	General Description of D-C Digital Voltmeters.....	Stansbury . . .	465
--------	--	-----------------	-----

Communication Division

61-827	Features of an Electronic Crosspoint PABX.....	Van Bosse, DeCicco . . .	471
61-722	Centrex Service: New Design for Group Telephone Service.....	Shea . . .	474
61-113	Pole-Zero Techniques Applied to V-F Telephone Lines.....	Fleming . . .	482
61-181	Noise and Intermodulation in Closed-Circuit TV.....	Collins, Williams . . .	486
61-825	Transmission Network of Electronic Crosspoint PABX.....	Kowalik . . .	491
61-826	Logical Control of Electronic Crosspoint PABX.....	Sanders . . .	496
61-822	Reliable Data Transmission Through Noisy Media.....	Melas . . .	501

Science and Electronics Division

61-732	Large-Scale On-Line Data Processing Systems.....	Levine . . .	505
60-1007	Computer Program for Preparing Wiring Diagrams..	Kirby, Rosenthal . . .	509
61-150	Magnetically Regulated D-C to D-C Converter Power Supply.....	Sager . . .	513
61-717	Regulated Power Supply Using Variable Output Oscillator.....	Jackson . . .	518
61-814	Forcing Circuitry: Sequential Building Blocks for Logic.....	Meade . . .	522
61-718	Inverter with Improved Commutation.....	McMurray, Shattuck . . .	531
61-725	A Comparison of Computers.....	Curl . . .	542
61-724	Results of Simulation Comparison of Control Computers.....	Sendzuk . . .	547

Education Committee

61-733	Educating Electrical Engineers for Professional Careers.....	Linke . . .	551
61-735	Motivation Through Challenge.....	Johnson, Clement . . .	554

(See inside back cover)

*Note to Librarians.* The six bimonthly issues of "Applications and Industry," March 1961-January 1962, will also be available in a single volume (no. 80) entitled "AIEE Transactions—Part II. Applications and Industry," which includes all technical papers on that subject presented during 1961. Bibliographic references to Applications and Industry and to Part II of the Transactions are therefore equivalent.

*Applications and Industry.* Published bimonthly by the American Institute of Electrical Engineers, from 20th and Northampton Streets, Easton, Pa. AIEE Headquarters: 345 East 47th Street, New York 17, N. Y. Address changes must be received at AIEE Headquarters by the first of the month to be effective with the succeeding issue. Copies undelivered because of incorrect address cannot be replaced without charge. Editorial and Advertising offices: 345 East 47th Street, New York 17, N. Y. Nonmember subscription \$8.00 per year (plus 75 cents extra for foreign postage payable in advance in New York exchange). Member subscriptions: one subscription at \$5.00 per year to any one of three divisional publications: Communication and Electronics, Applications and Industry, or Power Apparatus and Systems; additional annual subscriptions \$8.00 each. Single copies when available \$1.50 each. Second-class mail privileges authorized at Easton, Pa. This publication is authorized to be mailed at the special rates of postage prescribed by Section 132.122.

The American Institute of Electrical Engineers assumes no responsibility for the statements and opinions advanced by contributors to its publications.

Printed in United States of America

Number of copies of this issue 5,100



# The Effect of Variable Plasma Conductivity on MHD Energy Converter Performance

W. B. COE  
NONMEMBER AIEE

C. L. EISEN  
NONMEMBER AIEE

THE CONCEPT OF a magnetohydrodynamic (MHD) energy converter presents a means of converting some of the kinetic and thermal energy of a plasma directly into useful electrical work. The principle upon which its operation is based is that of passing a moving plasma through a transverse magnetic field, thereby inducing an electromotive force (emf) in the plasma. By completing the circuit through an external load, current is then permitted to flow. The analysis of such a device is complicated by the fact that the working fluid is a compressible medium, subject not only to the laws of compressible fluid flow, but to the laws of electromagnetism as well.

The potentialities of the MHD energy converter were investigated by Neuringer,<sup>1</sup> who applied the techniques of the calculus of variations to obtain an appropriate set of differential equations governing the conditions under which the maximum amount of useful electrical power could be generated. Although the equations he derived are quite general, they can be integrated in closed form only for special cases. One of the cases to which Neuringer applied these equations and obtained solutions was that of a converter with constant cross-sectional area and with constant plasma electrical conductivity. The results of this study showed that appreciable amounts of energy could be converted into useful electrical work.

A subsequent investigation undertook to obtain solutions to a set of differential equations governing the performance of the MHD converter. The equations were written for a varying channel cross-

sectional area and allowed the electrical conductivity of the plasma to depend on its local thermodynamic properties. These equations were programmed for numerical solution on an International Business Machine Corporation (IBM) 704 digital computer. The results of the computation permit the investigation of the performance of the converter for other than maximum energy conversion. A preliminary study<sup>2</sup> confined itself to a constant cross-sectional area converter with an assumed constant plasma conductivity.

This paper reports the results of a similar performance study of a constant area channel in which the electrical conductivity of the plasma is dependent on its local thermodynamic properties.

## Assumptions

There are a number of assumptions and approximations in the derivation of the equations governing the performance of the MHD energy converter. These include:

1. One-dimensional flow, i.e., the fluid dynamical state variables vary in the flow direction only, and not over the cross section.
2. Small magnetic Reynold's number, i.e., any effects on the fluid flow of second-

ary magnetic fields resulting from the induced current distribution are negligible, either because the secondary magnetic fields are small, or in the wrong direction to produce appreciable effects.

3. Negligible Hall current, i.e., the induced currents in the flow direction are small compared to those induced in the direction perpendicular to both the flow and to the magnetic fields.
4. The plasma behaves like a perfect gas.
5. The specific heats are constant.
6. The plasma is inviscid and non-heat-conducting so that the only dissipation is electrical.
7. The plasma is electrically neutral, i.e., no space charge sheaths are developed near the conducting walls.
8. Shock-free supersonic flow.
9. The channel walls are perfect electrical conductors.

## Equations

Let the flow of the plasma through the channel be in the  $x$  direction, as shown in the schematic diagram of the converter, Fig. 1. The channel cross-sectional area is given by the product  $y(x)z(x)$ , where  $y(x)$  is the channel height (or electrode spacing) and  $z(x)$  is the channel depth, both being prescribed functions of  $x$ , subject to the conditions that  $dy/dx$  and  $dz/dx$  are everywhere sufficiently small as required by the one-dimensional flow approximation. The applied magnetic field,  $B(x)$ , is in the  $z$  direction and is also a prescribed function of  $x$ . (Although the results discussed in the next section are for constant channel cross-section and constant magnetic-field intensity, both will here be considered as varying functions.) With the given directions for the magnetic field and the flow, an electric field is induced in the  $y$  direction, so that a potential difference thereby appears across the electrodes. The electrodes then are the terminals of the converter, and a load

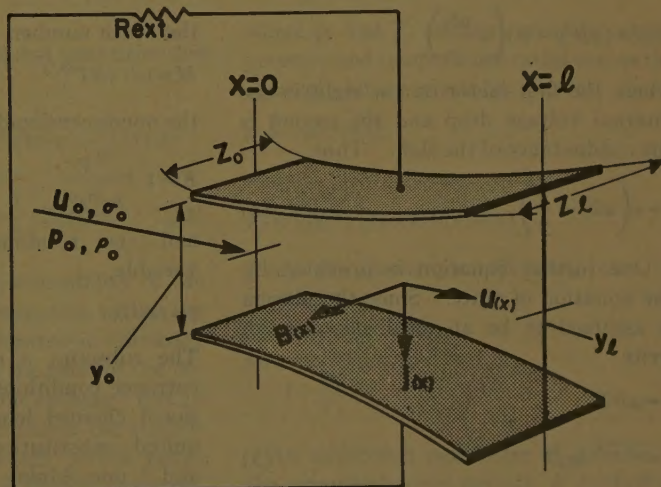


Fig. 1. Schematic of energy converter

Paper 60-1028, recommended by the AIEE Aerospace Transportation Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE Pacific General Meeting, San Diego, Calif., August 8-12, 1960. Manuscript submitted May 12, 1960; made available for printing April 17, 1961.

W. B. COE and C. L. EISEN are with the Republic Aviation Corporation, Farmingdale, N.Y.

This work was sponsored by the United States Air Force Office of Scientific Research under contract no. AF 49(638)-552 jointly with the Republic Aviation Corporation.



resistor connected across them will permit a current to flow.

The flow equations which apply to the plasma are the usual one-dimensional steady-state fluid dynamic equations of motion, with additional terms to account for the effect of the magnetic field on the flow. These equations are the continuity equation,

$$\frac{d}{dx}(\rho u y z)=0 \tag{1}$$

the momentum equation,

$$\rho u \frac{du}{dx} + \frac{dp}{dx} + jB = 0 \tag{2}$$

and the energy equation,

$$\frac{d}{dx} \left[ \rho u y z \left( c_p T + \frac{1}{2} u^2 \right) \right] + jB u y z - \frac{j^2 y z}{\sigma} = 0, \tag{3}$$

where  $\rho$  is the fluid density,  $u$  the fluid velocity,  $p$  the static pressure,  $j$  the induced current density,  $c_p$  the constant pressure specific heat,  $T$  the static temperature, and  $\sigma$  the local electrical conductivity of the plasma. In the momentum equation, the term  $jB$  is the reactive force per unit volume exerted by the magnetic field on the plasma. In the energy equation, the term  $jB u y z$  is the rate at which work is done per unit length of channel by the plasma against the reactive force, and would be the power delivered to the external load in the absence of internal losses. The term  $j^2 y z / \sigma$  is the power per unit length of channel dissipated within the plasma due to Joulean heating. For infinite conductivity, this term vanishes.

In addition to the fluid dynamic equations, an expression for the current density can be written for an assumed load on the converter. Let the load be such that the terminal voltage is  $V$ . Then the current in the slab of plasma of thickness  $dx$ , height  $y$ , and depth  $z$  is

$$jzdx = (uB y - V) \left( \sigma \frac{zdx}{y} \right)$$

where the first factor on the right is the internal voltage drop and the second is the conductance of the slab. Thus

$$i = \sigma \left( uB - \frac{V}{y} \right) \tag{4}$$

One further equation is provided by the equation of state. Since the plasma is assumed to be an ideal gas, one can write

$$p = \rho R T$$

$$= \frac{\gamma - 1}{\gamma} \rho c_p T \tag{5}$$

where  $R$  is the gas constant for the plasma, and  $\gamma$  is the specific heat ratio.

Equations 1-5 constitute five equations in the five unknown functions  $\rho$ ,  $u$ ,  $p$ ,  $T$ , and  $j$ . An immediate reduction to two equations in  $p$  and  $u$  can be made by using equations 4 and 5 and the integral of equation 1 to eliminate  $\rho$ ,  $T$ , and  $j$  in equations 2 and 3. One obtains then

$$\frac{du}{dx} + \frac{\gamma z}{m} \frac{dp}{dx} + \frac{B \sigma y z}{m} \left( uB - \frac{V}{y} \right) = 0 \tag{6}$$

$$\frac{d}{dx} \left( \frac{\gamma}{\gamma - 1} \gamma z u p + \frac{1}{2} m u^2 \right) + V \sigma z \left( uB - \frac{V}{y} \right) = 0 \tag{7}$$

where  $m$ , the mass flow rate, is the constant of integration of equation 1.

In order to make the parametric studies discussed later, it is convenient to nondimensionalize equations 6 and 7. For this purpose the following nondimensional quantities are introduced:

the nondimensional channel length,

$$\delta = \frac{B_0^2 y_0 z_0 \sigma_0}{m} l \tag{8}$$

the nondimensional velocity,

$$U = u / u_0$$

the nondimensional static pressure,

$$P = p / \rho_0 u_0^2$$

the nondimensional channel height,

$$Y = y / y_0$$

the nondimensional channel depth,

$$Z = z / z_0$$

the nondimensional electrical conductivity,

$$\bar{\sigma} = \sigma / \sigma_0 \tag{9}$$

the nondimensional magnetic field strength,

$$\bar{B} = B / B_0$$

the Mach number,

$$M = u / (\gamma R T)^{1/2}$$

the nondimensional load parameter,

$$K = 1 - \frac{V}{u_0 B_0 y_0}$$

and the nondimensional independent variable,

$$\xi = (\delta / l) x$$

The subscript  $o$  refers to the channel entrance conditions, and  $l$  is the dimensional channel length. Making the required substitutions into equations 6 and 7, one obtains

$$\frac{dU}{d\xi} + YZ \frac{dP}{d\xi} + \bar{B} \bar{\sigma} YZ \left[ \bar{B} U - \frac{1}{Y} (1 - K) \right] = 0 \tag{10}$$

and

$$\frac{d}{d\xi} \left( \frac{\gamma}{\gamma - 1} YZ U P + \frac{1}{2} U^2 \right) + (1 - K) \bar{\sigma} Z \left[ \bar{B} U - \frac{1}{Y} (1 - K) \right] = 0 \tag{11}$$

In equations 10 and 11,  $Y$  and  $Z$  are prescribed functions of  $\xi$ ,  $\bar{\sigma}$  is considered a known function of  $Y$ ,  $Z$ ,  $U$ , and  $K$  is the load parameter which can assume prescribed values between zero (open external circuit) and unity (short circuited external circuit). The determination of  $U(\xi)$  and  $P(\xi)$  over the interval  $0 \leq \xi \leq \delta$  constitutes the solution of equations 10 and 11.

To obtain an expression for the electrical conductivity it is assumed that the entire current is carried by the electrons. Thus

$$\sigma = \frac{j}{E} = \frac{n_e e u_e}{E} \tag{12}$$

where  $n_e$  is the number density of electrons,  $e$  the electronic charge,  $u_e$  the drift velocity of electrons through the gas, and  $E$  the electric field intensity. In general the electron mobility,  $u_e / E$ , is limited by collisions with neutral gas atoms and by long-range Coulomb interactions with other charged particles. The former process predominates in a weakly ionized gas, and the latter in a strongly ionized gas. If  $1 / \sigma_1$  is the resistivity of the gas due to electron collisions with neutral gas atoms and  $1 / \sigma_2$  the resistivity due to long-range Coulomb interactions, it is assumed then, that for a gas of intermediate ionization, in which both processes are significant, the resistivity is given by the sum of the resistivities due to the two processes:<sup>3</sup>

$$\frac{1}{\sigma} = \frac{1}{\sigma_1} + \frac{1}{\sigma_2} \tag{13}$$

For a weakly ionized gas in which the electric field is not large enough to cause ionizing collisions, the collision encounters between electrons and neutral gas atoms are mostly elastic. In such a gas Chapman's theory<sup>4</sup> for electron mobility gives

$$\frac{u_e}{E} = \frac{2}{3} \sqrt{\frac{2}{\pi}} \frac{e}{n Q} \left( \frac{1}{m_e k T} \right)^{1/2}$$

$$= \frac{2}{3} \sqrt{\frac{2}{\pi}} \frac{e}{Q p} \left( \frac{k T}{m_e} \right)^{1/2} \tag{14}$$

where  $n$  is the number density of gas atoms,  $Q$  the electron-gas atom collision cross section,  $k$  Boltzmann's constant,  $T$  the gas temperature, and  $p$  the pressure.



Although  $Q$  is a function of temperature, for this study it is taken as a suitable constant for the temperature ranges of interest.)

The number density of electrons as required by equation 12 can be obtained from Saha's equation,<sup>5</sup> assuming that the gas as a whole is electrically neutral and that the electrons are in equilibrium with the gas. It is further assumed that for temperatures of interest one can ignore all internal states of the gas atoms and their ions except the states of lowest energy. Saha's equation can then be written

$$n_e = n \frac{g_e g_i}{g_a} \left( \frac{2\pi m_e k T}{h^2} \right)^{3/2} \exp(-\varphi/kT) \quad (15)$$

$$= p \frac{g_e g_i}{g_a} \left( \frac{2\pi m_e k^{1/4} T^{1/4}}{h^2} \right)^{3/2} \exp(-\varphi/kT)$$

where  $g_e (=2)$  is the multiplicity of the electron,  $g_i$  and  $g_a$  are the respective multiplicities of the ground states of the ion and neutral atom,  $h$  is Planck's constant, and  $\varphi$  is the ionization potential.

Substituting into equation 12 for the mobility as given by equation 14, and for the electron particle density as given by equation 15, one has for a weakly ionized gas

$$\eta_i(p, T) = \frac{(2^9 \pi)^{1/4} e^2 m_e^{1/4} k^{3/4}}{3 h^{3/2}} \times \left( \frac{g_e g_i}{g_a} \right)^{1/2} \frac{1}{Q} \frac{T^{3/4}}{p^{1/2}} \exp(-\varphi/2kT) \quad (16)$$

For a fully ionized gas, the electrical conductivity, as given by Spitzer and Härm,<sup>6</sup> is, in mks (meter-kilogram-second) units,

$$\sigma = \frac{2^{13/2} \pi^{1/2} \kappa_0^2 \gamma_E k^{1/2}}{m_e^{1/2} Z e^2} \cdot \frac{T^{3/2}}{\ln(qc^2)} \quad (17)$$

where  $\kappa_0$  is the permittivity of free space,  $\gamma_E$  a correction factor to account for electron-electron encounters, and  $Ze$  the ionic charge:

$$= \frac{4\pi \kappa_0^{3/4} m_e (kT)^{1/2}}{Z(Z+1)^{1/2} e^2 n_e^{1/2}} \quad (18)$$

$$= \left( \frac{3kT}{m_e} \right)^{1/2} \quad (19)$$

For a singly ionized gas,  $Z=1$  and  $\gamma_E = 1.5816$ .

For a fully ionized gas,

$$\sigma = \frac{p}{2kT} \quad (20)$$

Equations 17-20 combine to express  $\sigma_2$  as a function of  $p$  and  $T$ . The results together with equation 16 are substituted into equation 13 to provide the conduc-

tivity function, which is then nondimensionalized according to equation 9, where  $\sigma_0$  is the conductivity evaluated at the entrance to the channel. Equations 10 and 11, however, require that the conductivity be expressed as a function of the nondimensional variables. Thus the substitutions

$$p = \gamma M_o^2 p_o P$$

and

$$T = \gamma M_o^2 T_o YZUP$$

must be made in  $\bar{\sigma}$ .

The power output of the converter is given by

$$\dot{P} = V \int_0^l j_z dx = V \int_0^l \sigma z \left( uB - \frac{V}{y} \right) dx$$

where the current density is obtained from equation 4. Substituting for the integrand from equation 7, the output power can be expressed in the following convenient form:

$$\dot{P} = h_T(0) - h_T(l)$$

where

$$h_T(x) = \frac{\gamma}{\gamma-1} y z u p + \frac{1}{2} m u^2$$

is the stagnation enthalpy flux at  $x$ . A measure of the output power is therefore provided by the conversion effectiveness,  $\eta_c$ , which is defined as the ratio of the output power to the stagnation enthalpy flux at the entrance to the channel:

$$\eta_c = \frac{\dot{P}}{h_T(0)} = 1 - \frac{h_T(l)}{h_T(0)} = 1 - \frac{T_T(l)}{T_T(0)} \quad (21)$$

where  $T_T$  is the stagnation temperature,  $c_P$  being assumed constant. The nondimensional stagnation enthalpy flux is defined by

$$H_T(\xi) = \frac{h_T(x)}{m u_o^2} \quad (22)$$

In terms of nondimensional quantities this can be written

$$H_T(\xi) = \frac{\gamma}{\gamma-1} Y(\xi) Z(\xi) U(\xi) P(\xi) + \frac{1}{2} [U(\xi)]^2 = \frac{1}{\gamma-1} \left[ \frac{U(\xi)}{M(\xi)} \right]^2 + \frac{1}{2} [U(\xi)]^2 \quad (23)$$

A measure of the irreversibility of the flow due to Joulean dissipation within the fluid is given by the isentropic efficiency  $\eta_s$ , which is defined as the ratio of the decrease in the stagnation enthalpy flux in the channel to the decrease in the stagnation enthalpy flux which would result from an isentropic flow process terminat-

ing in an exit stagnation pressure equal to that of the actual flow. Thus

$$\eta_s = \frac{h_T(0) - h_T(l)}{h_T(0) - h'_T(l)} = \frac{\eta_c}{1 - \frac{h'_T(l)}{h_T(0)}} = \frac{\eta_c}{1 - \frac{T'_T(l)}{T_T(0)}}$$

again assuming the constancy of the specific heat. Here  $h'_T(l)$  and  $T'_T(l)$  are respectively the exit stagnation enthalpy flux and temperature for the isentropic flow process. One can also write for the isentropic process

$$\frac{T'_T(l)}{T_T(0)} = \left[ \frac{p_T(l)}{p_T(0)} \right]^{(\gamma-1)/\gamma}$$

where  $p_T$  is the total pressure, and where  $p_T(l)$  is the exit total pressure for both the actual flow and the isentropic flow. Thus

$$\eta_s = \frac{\eta_c}{1 - \left[ \frac{p_T(l)}{p_T(0)} \right]^{(\gamma-1)/\gamma}}$$

To compute  $\eta_s$ , therefore, one must express the total pressure ratio,  $p_T(l)/p_T(0)$ , across the channel in terms of the solutions of equations 10 and 11. To do this one notes that the stagnation to static pressure and temperature are given respectively by

$$\frac{p_T}{p} = \left[ 1 + \frac{1}{2} (\gamma-1) M^2 \right]^{\gamma/(\gamma-1)} \quad (24)$$

$$\frac{T_T}{T} = 1 + \frac{1}{2} (\gamma-1) M^2$$

so that

$$\frac{p_T(l)/p_T(0)}{p(l)/p(0)} = \left[ \frac{T_T(l)/T_T(0)}{T(l)/T(0)} \right]^{\gamma/(\gamma-1)}$$

or

$$\left[ \frac{p_T(l)}{p_T(0)} \right]^{(\gamma-1)/\gamma} = \frac{p_r^{(\gamma-1)/\gamma}}{T_r} \times \frac{T_T(l)}{T_T(0)} = \frac{p_r^{(\gamma-1)/\gamma}}{T_r} (1 - \eta_c)$$

where  $p_r$  and  $T_r$  are the respective static pressure and temperature ratios across the channel and are given by

$$p_r = \gamma M_o^2 P(\delta)$$

and

$$T_r = \gamma M_o^2 Y(\delta) Z(\delta) U(\delta) P(\delta)$$

Thus the isentropic efficiency can be written

$$\eta_s = \frac{\eta_c}{1 - \frac{p_r^{(\gamma-1)/\gamma}}{T_r} (1 - \eta_c)}$$

An additional parameter of interest is the generated power density,  $\dot{p}$ , defined as



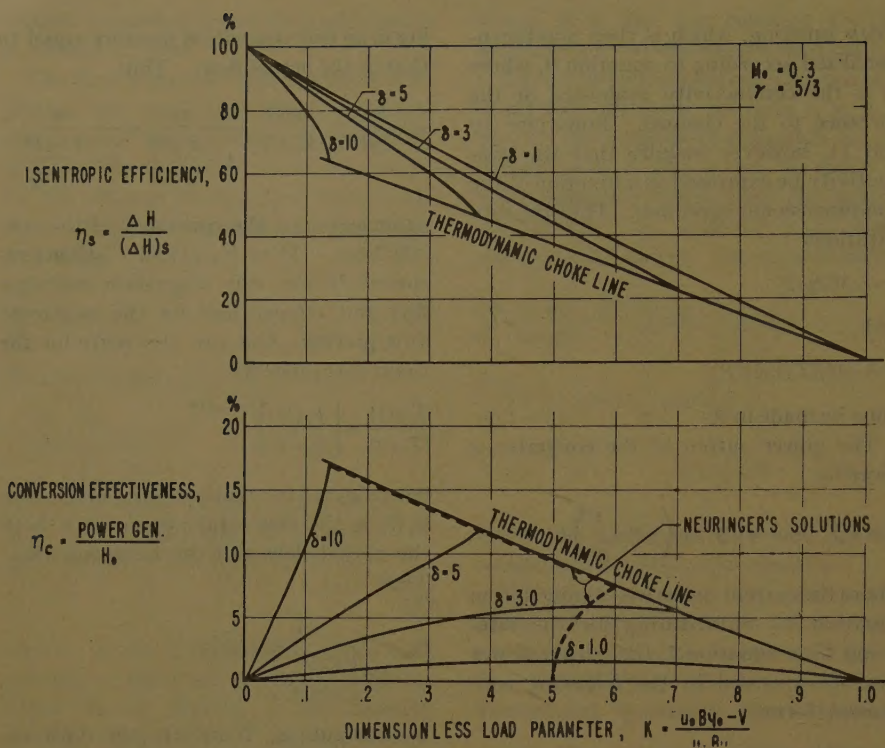


Fig. 2. Subsonic performance, constant area, constant conductivity

the ratio of the power output to the volume,  $v$ , of the channel:

$$\hat{p} = \frac{h_T(0)\eta_c}{v} \quad (25)$$

From equations 22 and 23, it is seen that the entrance enthalpy flux can be written

$$h_T(0) = mu_0^2 \left[ \frac{1}{(\gamma-1)M_0^2} + \frac{1}{2} \right] \quad (26)$$

The initial velocity may be expressed as

$$u_0 = (\gamma RT_0)^{1/2} M_0 \quad (27)$$

and by equation 8, the mass flow ratio may be given by

$$m = \frac{B_0^2 y_0 z_0 \sigma_0 l}{\delta} \quad (28)$$

For the volume of the channel, one has

$$\begin{aligned} v &= \int_0^l yz dx \\ &= \frac{y_0 z_0 l}{\delta} \int_0^\delta YZ d\xi \\ &= \frac{y_0 z_0 l}{\delta} \bar{v} \end{aligned} \quad (29)$$

where a nondimensional volume is defined by

$$\bar{v} = \int_0^\delta YZ d\xi \quad (30)$$

Substituting the expressions given by equations 26-29 into equation 25, one obtains

$$\begin{aligned} \hat{p} &= \frac{\gamma}{\gamma-1} B_0^2 \sigma_0 RT_0 \left[ 1 + \frac{1}{2} (\gamma-1) M_0^2 \right] \frac{\eta_c}{\bar{v}} \\ &= \frac{\gamma}{\gamma-1} B_0^2 \sigma_0 RT_T(0) \frac{\eta_c}{\bar{v}} \end{aligned} \quad (31)$$

Except for the channel shape, the mag-

netic field distribution, and the gas data ( $\gamma$ ,  $R$ ,  $g_i$ ,  $g_a$ ,  $\varphi$ , and  $Q$ ), the input parameters needed to compute  $\eta_c$  and  $\eta_s$  are  $M_0$ ,  $\delta$ ,  $K$ ,  $T_0$ , [or  $T_T(0)$ ], and  $p_0$  [or  $p_T(0)$ ], the temperature and pressure being needed only in the conductivity function. The computation of  $\hat{p}$  also requires  $B_0$ .

## Discussion and Results

The nondimensional equations 10 and 11 were solved numerically, using the IBM 704 digital computer. The performance of the MHD energy converter was calculated for constant channel cross-sectional area ( $Y=1$  and  $Z=1$ ) and for constant magnetic field ( $\bar{B}=1$ ). The specific heat ratio  $\gamma$  was taken as 5/3, representing either a monatomic gas or fully ionized plasma.

The performance was computed for constant area energy converter; first with an assumed constant plasma conductivity ( $\sigma=1$ ), then, for a specified working fluid with the plasma conductivity, given by equations 9, 13, 16, and 17 as a function of the local thermodynamic properties of the plasma. The performance was obtained for a variety of subsonic and supersonic entrance Mach numbers  $M_0$ , and for a variety of nondimensional channel lengths  $\delta$ , the load parameter  $K$ , ranging from open circuit ( $K=0$ ) to short circuit ( $K=1$ ).

Figs. 2 and 3 present the performance

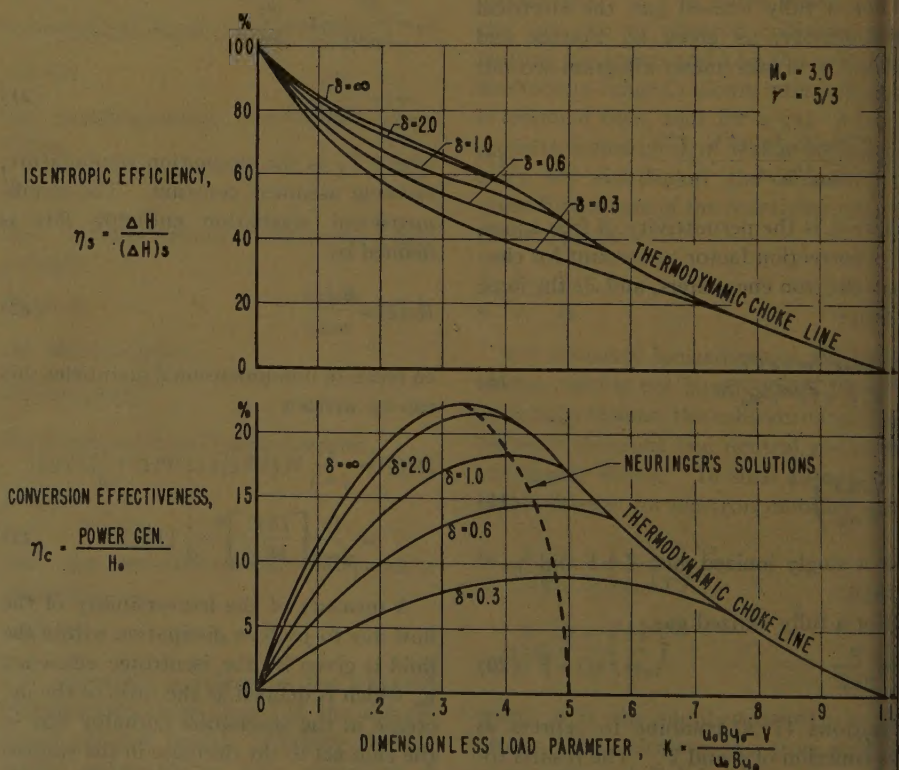


Fig. 3. Supersonic performance, constant area, constant conductivity



the constant area energy converter for assumed constant plasma conductivity.

The performance of a subsonic energy converter is shown in Fig. 2 for an entrance Mach number of 0.3. The isentropic efficiency  $\eta_s$  and the conversion effectiveness  $\eta_c$  are plotted as functions of the load parameter  $K$  for a range of channel lengths  $\delta$ . For small values of the conversion effectiveness increases and the isentropic efficiency decreases as  $K$  is increased from zero. At approximately  $K=1/2$  the conversion effectiveness is a maximum, and further increases in  $K$  cause  $\eta_c$  to decrease until at  $K=1$  the energy converted becomes zero. The isentropic efficiency continually decreases with increasing  $K$  until it also goes to zero at  $K=1$ . When  $K=1$ , the converter is short-circuited and all the gas dynamic energy removed from the plasma is returned in the form of Joulean dissipation. For larger values of the channel length,  $\delta$ , the load parameter,  $K$ , cannot be increased beyond a certain limiting value. This limit is imposed by the flow becoming sonic at the exit of the channel. When this occurs the flow is choked. The locus of this condition is indicated by the line labeled "thermodynamic choke line." It is thus seen that for shorter channels the maximum energy conversion occurs when the loading is chosen in a manner which best matches the energy output to the Joulean dissipation; for longer channels the maximum energy conversion is dictated by the thermodynamic choking of the flow.

At a fixed value of loading, increasing the length of the converter causes an increase in  $\eta_c$  and a decrease in  $\eta_s$ . It also increases the exit flow Mach number until the choke condition is reached. For the subsonic converter the flow velocity increases along the channel length. This in turn produces a larger induced emf and requires, therefore, larger internal voltage drops to match the terminal voltage specified by  $K$ . Each subsequent increment in length added to the channel converts additional energy, but does so with poorer local efficiency.

The performance for the supersonic converter is shown in Fig. 3 for an entrance Mach number of 3.0. The "thermodynamic choke line" again notes the operating conditions where the exit flow Mach number is unity. As compared to the subsonic converter, at a fixed value of loading, both the conversion effectiveness and the isentropic efficiency of the supersonic converter increase as the length of the channel is increased. For this case the flow velocity

decreases along the length of the channel, causing a reduction in the induced emf. This requires, therefore, a smaller internal voltage drop in order to match the specified terminal voltage. Thus subsequent increments of length added to the channel convert energy with a local efficiency that is better than the average for the channel. For lightly loaded operation, i.e., small values of  $K$ , the drop in velocity causes the induced emf to approach the terminal voltage asymptotically. When this occurs further lengthening of the channel produces no change in the flow conditions or the energy converted. The line denoted by  $\delta = \infty$  indicates this condition.

Figs. 2 and 3 also show that the maxima of the conversion effectiveness agree with the results of Neuringer's optimum power analysis.<sup>1</sup> An interesting characteristic is the value of  $K$  at which these maxima occur. If the plasma acted as a perfectly incompressible fluid, or a solid conductor, maximum conversion would occur when the energy delivered to the external load is equal to that lost by Joulean dissipation within the converter. For an incompressible fluid this condition is obtained when  $K=1/2$ , since the induced emf is constant down the channel. As seen in Figs. 2 and 3, for short channels both the subsonic and supersonic converter also require  $K \sim 1/2$  for maximum conversion. As the channel length of

both converters is increased the value of  $K$  at maximum conversion deviates from  $K=1/2$ , the subsonic converter requiring larger values of  $K$  in the unchoked region, and the supersonic converter always requiring smaller values of  $K$ . For these larger channel lengths the compressibility effects are significant and the local loading of each element of the channel is related to that of all the other elements through the gas dynamics involved. As a result, the condition for maximum conversion is no longer simply that required by a solid or incompressible conductor.

In order to study the effect of a variable conductivity plasma upon the performance of the MHD energy converter, a choice must be made of the working fluid and its initial conditions. For this investigation cesium vapor was chosen as the working fluid, the choice being dictated by its low value of ionization potential,  $\phi=3.87$  volts, which will provide reasonable values of plasma conductivity for temperatures in the thermal range. The performance was computed for two values of entrance stagnation temperature,  $T_{70}$ . The lower value,  $T_{70}=2,000$  K (degrees Kelvin), represents the lowest usable plasma temperature at which present-day materials can conceivably operate. The higher value of 5,000 K has been included to indicate the increased performance capability that could be obtained, without regard for the materials

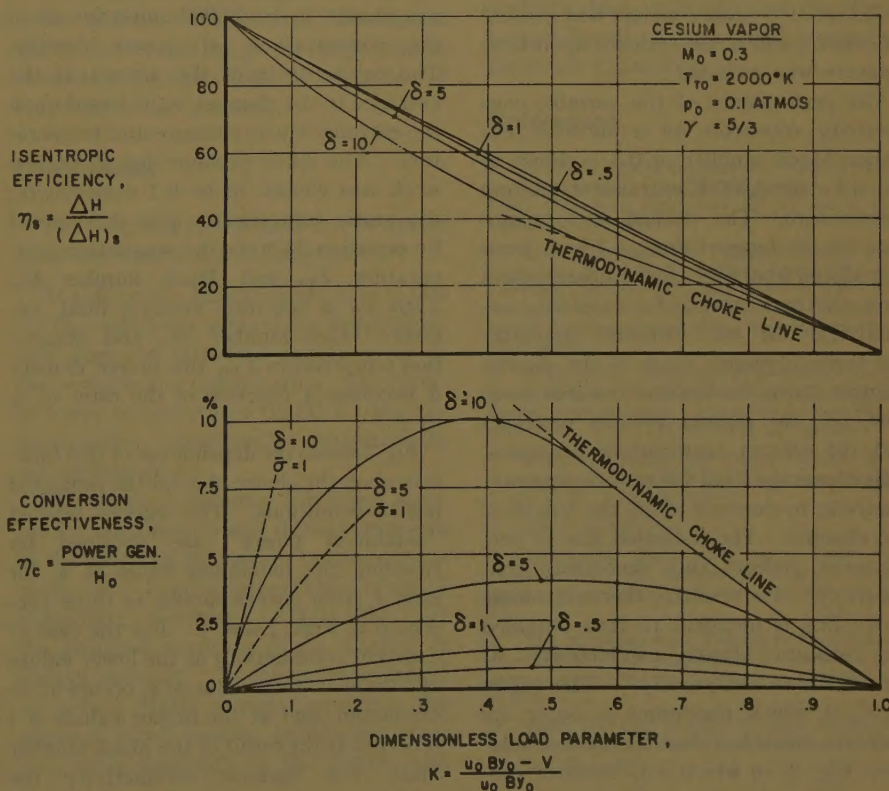


Fig. 4. Subsonic performance, constant area, variable conductivity



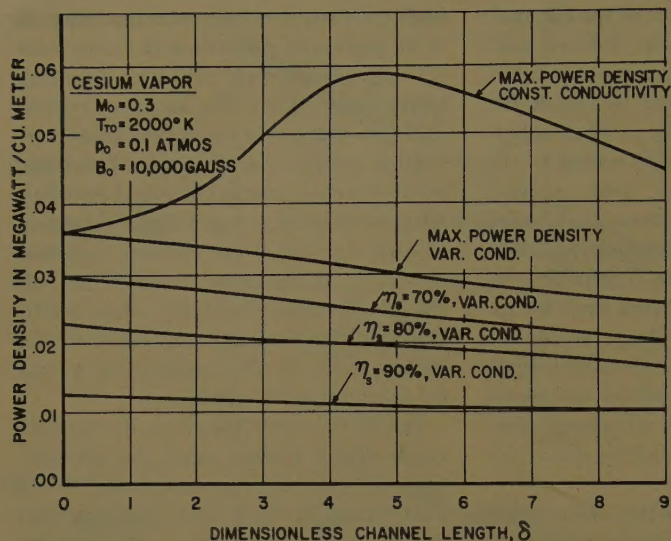


Fig. 5. Power density as a function of channel length

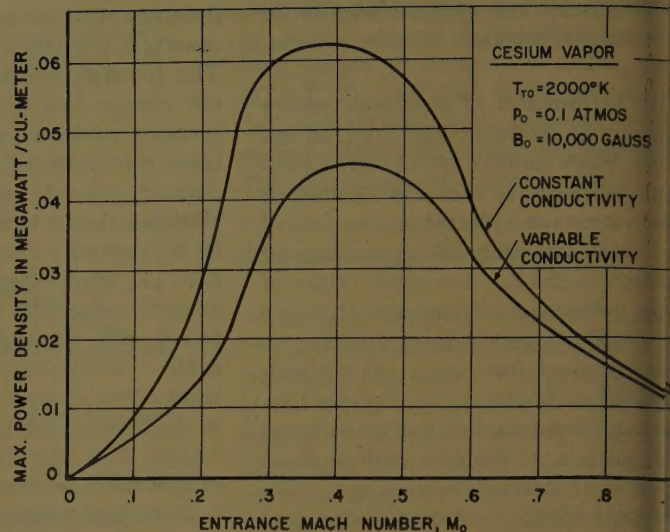


Fig. 6. Maximum power density for  $T_{To} = 2,000$  K

problem. All the performance computations reported here, for variable plasma conductivity, have been made for an entrance static pressure,  $p_o$ , of 0.1 atmosphere. This represents the lowest value for the present choice of working fluid and temperatures that will not violate the assumption of a scalar plasma conductivity.<sup>7</sup>

The conductivity function requires the collision cross section of electrons with atoms as a function of temperature. Because good data were not available in the required temperature range, a value of  $3.25 \times 10^{-18}$  square meters was chosen<sup>8</sup> for cesium, and no dependence upon temperature was included.

The performance of the variable conductivity converter for a subsonic entrance Mach number of 0.3 is shown in Fig. 4 for the 2,000 K entrance stagnation temperature. The dashed performance lines for the larger values of  $\delta$  have been reproduced from Fig. 2 for comparison and show the performance for constant conductivity. In the subsonic converter the thermodynamic state of the plasma changes down the channel towards sonic flow, i.e., the plasma velocity increases and the plasma temperature decreases. This causes the local value of plasma conductivity to decrease along the length of the channel. The indicated loss in performance reflects this decreasing conductivity. The resulting thermodynamic choke line is identical to that obtained for constant plasma conductivity, as can be shown analytically. The values of  $K$  at which maximum  $\eta_c$  occur are lower for variable  $\sigma$  than for the computation (Fig. 2) in which  $\sigma$  is assumed constant. For the case of supersonic entrance Mach numbers, the local value of

conductivity increases along the length of the channel and results in an increase in converter performance. The consequences of this characteristic will be indicated in some of the subsequent results.

For a constant area channel the non-dimensional volume, as given by equation 30, reduces to  $\bar{v} = \delta$ . The power density  $\hat{p}$  given by equation 31 then becomes

$$\hat{p} = \frac{\gamma}{\gamma - 1} B_o^2 \sigma_o R T_{To} \frac{\eta_c}{\delta} \quad (32)$$

The constant magnetic field strength,  $B_o$ , was chosen to be 10,000 gauss for all of the computations of power density. The conductivity of the plasma at the entrance to the channel,  $\sigma_o$ , is based upon the entrance static pressure and temperature. The static pressure for all of this work was chosen to be 0.1 atmosphere; the static temperature was determined by equation 24 from the stagnation temperature  $T_{To}$  and Mach number  $M_o$ . Thus for a specified working fluid, entrance Mach number  $M_o$  and stagnation temperature  $T_{To}$ , the power density  $\hat{p}$  becomes a function of the ratio of  $\eta_c$  to  $\delta$ .

Fig. 5 shows the dependence of this function upon the choice of  $\delta$ , for the indicated initial conditions. The curves labeled "maximum power" are obtained by choosing the maximum value of  $\eta_c$  for each  $\delta$ , from curves similar to those presented in Figs. 2 and 4. For the case of constant conductivity at the lower values of  $\delta$  the maximum value of  $\eta_c$  occurs as an extremum, and at the higher values of  $\delta$  it occurs as the result of the Mach number limit. For variable conductivity the performance is substantially reduced, and at all  $\delta$ 's calculated,  $\eta_c$  occurs as an

extremum. Also included in Fig. 5 are power densities for various constant values of isentropic efficiency,  $\eta_s$ . These curves are included to show how the power density would be decreased from its maximum if high isentropic efficiency were an important design consideration.

For temperatures of the order of 2,000 K the conductivity of a cesium plasma varies predominantly as an exponential function of temperature as indicated in equation 16. The entrance static temperature decreases for increasing values of entrance Mach numbers and fixed values of  $T_{To}$ . Although this results in a large reduction in  $\sigma_o$ , the maximum power density for small subsonic Mach numbers, increases with Mach number, as shown in Fig. 6 for  $T_{To} = 2,000$  K. As seen from equation 32 this is possible only because the ratio of  $\eta_c$  to  $\delta$  increases faster than  $\delta$  decreases. As the entrance Mach number is further increased a point is reached after which the maximum power density decreases, since at the high Mach number, the decrease in  $\sigma_o$  is the dominant factor in equation 32. The value of power density for supersonic entrance Mach numbers has not been calculated for  $T_{To} = 2,000$  K due to the extremely low values of plasma conductivity. The sensitive nature of the conductivity function at temperatures in this range causes an appreciable variation in the power density that can be achieved when the plasma conductivity is allowed to vary in a realistic manner as opposed to holding it a constant. The performance shown in Fig. 6 indicates an overestimation of about 40% in power density if the conductivity is assumed to be constant.

The performance shown in Fig. 7 is for an entrance stagnation temperature



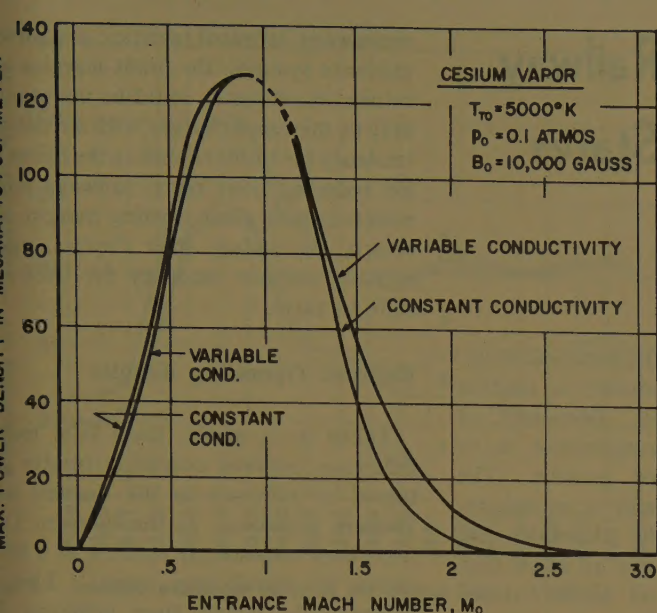


Fig. 7. Maximum power density for  $T_{T_0} = 5,000$  K

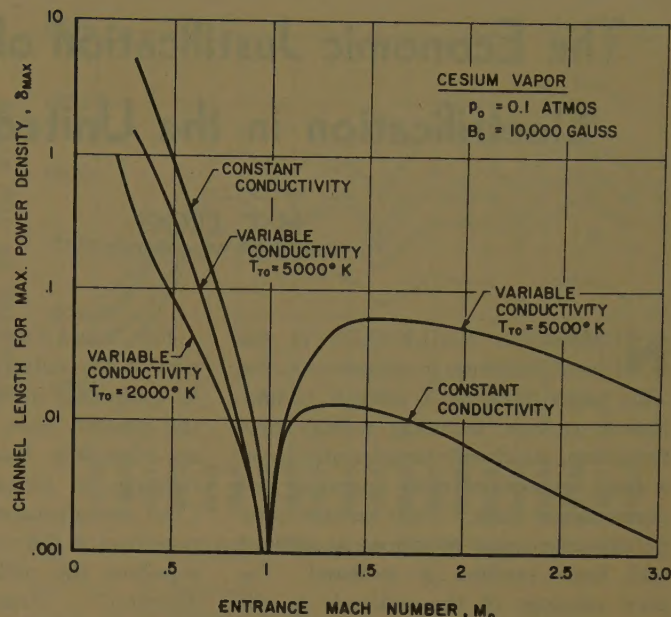


Fig. 8. Converter length for maximum power density

5,000 K. The large increase in the level of power density is mainly due to the exponential increase in the level of conductivity, although the increased stagnation temperature itself causes an increase by a factor of 2.5. This indicates that for small increases in the stagnation temperature above 2,000 K large gains in power density could be achieved. For a stagnation temperature of 5,000 K the variation of  $\sigma$  with changes in temperature is much smaller than at the lower temperature. This is indicated in Fig. 7 by the small difference between the performance with and without variable conductivity. Also because of this smaller variation in the performance computation may be carried into the supersonic range. The maximum value of power density now occurs at a higher value of entrance Mach number.

The length of channel,  $\delta$ , required to achieve the maximum power densities shown in Figs. 6 and 7 is plotted in Fig. 8 as a function of entrance Mach number. For constant plasma conductivity the required  $\delta$  is independent of temperature. The effect of variable conductivity is to shorten the channel length for subsonic Mach number and lengthen the channel for supersonic Mach number.

With a specified value of 2,000 K for the entrance stagnation temperature and an entrance static pressure of 0.1 atmos where the maximum power density, for a variable plasma conductivity, occurs at an entrance Mach number of  $M_0 \sim 0.45$ , as shown in Fig. 6. The nondimensional channel length required for maxi-

mum power density is  $\delta \sim 0.1$ , as given by the length functions of Fig. 8. From the definition of  $\delta$  (equation 8) one can write

$$l = \frac{\rho_0 u_0 \delta}{B_0^2 \sigma_0} = \frac{M_0 \delta p_0 \gamma^{1/2}}{\sigma_0 B_0^2 (RT_0)^{1/2}} \quad (33)$$

where  $l$  is the dimensional length of the channel. Substituting into equation 33 for the case cited above, one obtains  $l \sim 50$  centimeters. For a channel height and width equal to its length the output power will be  $\sim 5$  kilowatts. The size of the converter and the power output can be controlled by the choice of the value of entrance static pressure.

## Conclusions

The following conclusions can be drawn from the results of this analysis:

At sufficiently high temperatures (5,000 K) the constant conductivity approximation does not introduce serious error in the performance computation. However, at lower temperatures (2,000 K) the approximation significantly overestimates the performance.

The required entrance flow Mach number for maximum power density increases as the entrance stagnation temperature is increased. For the temperature range of immediate interest (2,000–5,000 K) the required entrance Mach number is always subsonic. In the construction of an actual MHD energy converter, the subsonic Mach number requirement permits more flexibility in the design, and

lends more validity to the approximations which are made in the analysis than would be true for a supersonic energy converter.

The attainable power density increases sharply with increases in entrance stagnation temperature. Since the maximum temperature that an MHD energy converter can operate at is limited by the capabilities of the materials involved, the reward of large increases in performance for small increases in operating temperatures provides a strong incentive for sophistication in design and for the accomplishment of even small increases in materials capability.

## References

1. OPTIMUM POWER GENERATION USING A PLASMA AS A WORKING FLUID, J. L. Neuringer. *Journal of Fluid Mechanics*, New York, N. Y., vol. 7, pt. 2, 1960, pp. 287–301.
2. THE PERFORMANCE CHARACTERISTICS OF A CONSTANT AREA MHD ENERGY CONVERTER WITH ASSUMED CONSTANT PLASMA CONDUCTIVITY, C. L. Eisen, W. B. Coe. *Report no. PPL-TN-60-7*, Republic Aviation Corporation, Farmingdale, N. Y., Mar. 1960.
3. ELECTRICAL CONDUCTIVITY OF HIGHLY IONIZED ARGON PRODUCED BY SHOCK WAVES, S. Lin, E. L. Resler, A. Kantrowitz. *Journal of Applied Physics*, New York, N. Y., vol. 26, 1955, pp. 95–109.
4. KINETIC THEORY OF GASES (book), E. H. Kennard. McGraw-Hill Book Company, Inc., New York, N. Y., 1938, pp. 470–73.
5. *Ibid.*, pp. 243–45.
6. TRANSPORT PHENOMENA IN A COMPLETELY IONIZED GAS, L. Spitzer, R. Härm. *Physical Review*, New York, N. Y., vol. 89, 1953, pp. 977–81.
7. THE MATHEMATICAL THEORY OF NON-UNIFORM GASES (book), S. Chapman, T. G. Cowling. Cambridge University Press, Cambridge, England, 1952, p. 328.
8. BASIC DATA OF PLASMA PHYSICS (book), Sanborn C. Brown. John Wiley & Sons, Inc., New York, N. Y., 1959, p. 6.



# The Economic Justification of Railway Electrification in the United States

H. C. CROSS  
ASSOCIATE MEMBER AIEE

**N**UMEROUS RAILROADS in the United States are in economic distress. This paper presents a modern technological package, including railway electrification, which will permanently lower a large segment of their operating and maintenance costs. Their current competitive economic situation is outlined and their position is reviewed. The poor earnings of the railroads in the northeast sector—ICC (Interstate Commerce Commission) Eastern District—relative to those in the other sectors—ICC Western District—are shown.

New types of communication and signal-interlocker systems, electric locomotives converted from existing diesel-electric locomotives, and the commercial-frequency high-voltage system of electrification are included in the technological package, which is described in detail, together with its advantages.

An economic study shows the annual economic results for one year before conversion and each of 12 years thereafter, using a typical Eastern District railroad, operating inland from the eastern seaboard, as its basis. For purposes of comparison, results are shown with both dieselization and electrification. The package with electrification produces major and permanent savings not offered with dieselization. Moreover, the net capital investment is less with electrification; major increases are made in freight train speeds, and equal or increased capacity is provided.

For railroads in the Eastern District which now have relatively heavy traffic, and for the heavy traffic routes which will develop from future railway mergers, this package is of great economic importance.

## Railroads as a Growth Industry

The railroad transportation industry in the United States is regarded as unprogressive. Freight traffic, as measured in ton-miles, remains more or less unchanged, but passenger traffic, measured in terms of passenger-miles, is continuously decreasing, even though, in over-all terms, both intercity freight and passenger traffic is steadily increasing.

Figs. 1 and 2 show the historical situation of the railway industry in relation to competing industries. Obviously, at the present time, transportation is in an effectively balanced position. The railways no longer hold a monopoly. "The manufacturer, the wholesaler, the retailer, and the customer all know that somehow the milk will arrive on the doorstep." From the viewpoint of national self-interest, the railway industry should share in the increasing traffic by, at least, maintaining its existing percentages of the total.

Fig. 3 shows the operating results for Class I railroads and the competing industries. The railways' profit margin and return on net assets have declined steadily since 1955. The 1959 margin was the same as for total manufacturing, twice that for motor carriers and air transport, and only slightly lower than that for interstate bus companies and inland water carriers. The return on net assets of competing companies ranges from 1.8 to 5.1 times that of the railways. This is basically the result of (1) overstatement of the railroad industry net worth growing out of conservative depreciation policies and (2) the capital investment in right of way, tracks, and signals, for which competing industries make no comparable investment. This situation is clearly brought out in the graph at the right-hand side of Fig. 3, showing net assets per dollar of sales. Trucks, buses, and airplanes operate at abnormally low net assets per dollar of sales because of indirect subsidies which, as national policy, were granted them in order to break the previously existing railway transportation monopoly.

If the competing industries were to pay the full costs of independent throughways,

Paper 61-208, recommended by the AIEE Land Transportation Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE Winter General Meeting, New York, N. Y., January 29-February 3, 1961. Manuscript submitted November 23, 1960; made available for printing March 14, 1961.

H. C. Cross is with Westinghouse Electric International Company, A Division of Westinghouse Electric Corporation, New York, N. Y.

The author acknowledges assistance rendered in various ways by H. F. Brown, L. W. Birch, T. A. Benner, E. K. Bloss, R. L. Kimball, C. R. Kingston, A. G. Oehler, and J. R. Shepard.

waterways, terminal facilities, and airway guidance systems, the profit margins and return on net assets could be maintained only by increased charges, with a resultant tendency for traffic to shift to the railways. By reducing their costs, railways could maintain both their existing margin and return, and reduce their charges, again with a resultant tendency for traffic to shift to them.

## Railroad Operating Results

Taken as a whole, there is a marked difference between operating results obtained by railroads in the Eastern and Western Districts. In the Western District, they consistently operate at a level above the break-even point. This is substantially higher than achieved by Eastern District railways, even though their traffic in gross ton-miles per mile of route is about one third higher and their total revenue per mile of route is approximately 80% higher than in the Western District. Western's margin on sales during 1958 was about 8.3% as compared with Eastern's margin of approximately 2.5%.

Relatively, Eastern railways have greater trackage per mile of route, particularly yard switching tracks. The particularly high proportion of yard-switching tracks indicates relatively excessive terminal costs undoubtedly arising from duplication of facilities.

They encounter more severe competition from (1) trucks operating over highways provided at public expense and shipping on inland waterways and the St. Lawrence Seaway route into the midwest. As a whole, they normally operate uncomfortably close to the break-even position as to earnings. During a recession or even temporary major suspension such as was caused by the steel strike in 1959, individual roads immediately take a loss. In view of the low margin on sales, rate reductions to handle existing traffic must be limited. On a contributed margin—or out-of-pocket basis, rate reductions must either attract additional customers without disturbing the existing volume, or these reductions must be in such large volume as will more than offset any reduction of income from existing traffic.

The imperative need seems to be reduction in operating and maintenance costs attained with a minimum of capital expenditure. Also, to assist in attracting additional traffic to the railroads, greater increased speed is required of freight trains.



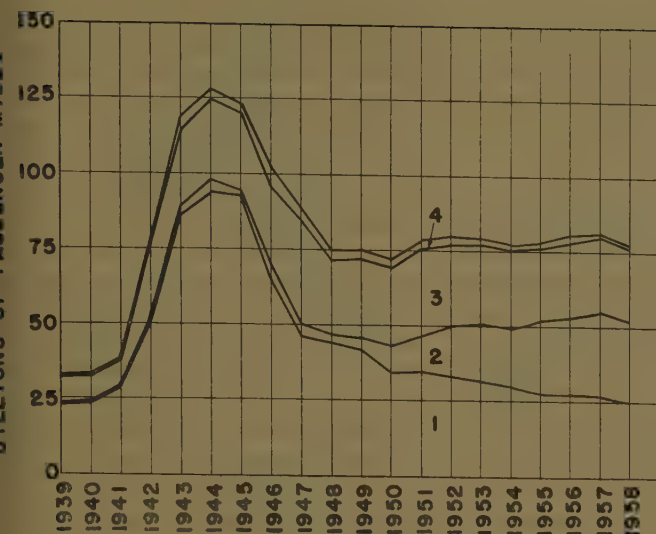


Fig. 1. Intercity passenger traffic in the United States, 1939-58

- 1—Railroads
  - 2—Airplanes
  - 3—Buses
  - 4—Water carriers and electric railways
- (Numbers refer to area between lines)

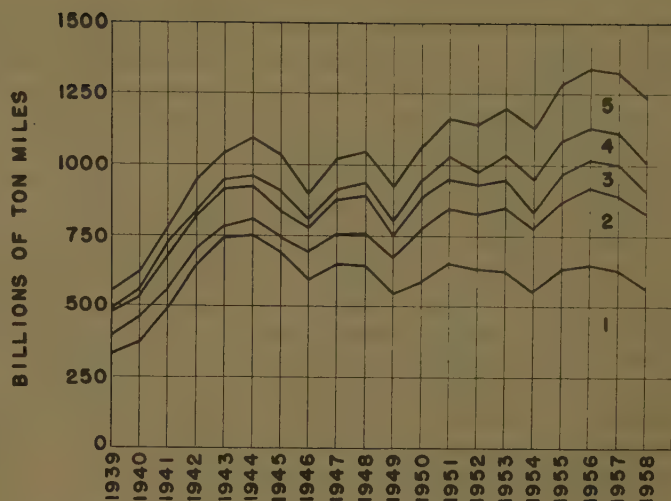


Fig. 2. Intercity freight traffic in the United States, 1939-58

- 1—Railroads
  - 2—Motor carriers
  - 3—Great Lakes
  - 4—Rivers and canals
  - 5—Oil pipelines
- (Numbers refer to areas between lines)

## Technological Package

### SUITABLE STANDARD SYSTEM FOR RAILWAY ELECTRIFICATION

A standard electrification system having a world-wide application has come to existence in recent years, making use of a commercial-frequency single-phase, 25,000-volt power supply to a catenary-suspended trolley over railway tracks. This system became practical with the development of rectifier-type locomotives in which the 25,000-volt power supply is stepped down to a voltage suitable for application to rectifiers supplying a pulsating direct current to standard series commutator traction motors.

The ignitron rectifier locomotive initially was developed in the United States and was adopted extensively by French manufacturers, both for use in their own country and for export. The British railways also have adopted this system for electrification of their trackage, and are actively engaged in installing the power supply system in various regions. Their national manufacturers are also engaged in production of necessary rectifier-type locomotives. The experience of the French National Railways with this system of electrification,<sup>1</sup> coupled with the theoretical studies which have been made in the United States,<sup>2,3</sup> serve to remove doubts as to this system's practicality.

Test results, covering unbalanced loadings of turbine generators, suggest that this system could have application gen-

erally in the United States, and particularly in regions where there exists a power-generating system of high capacity, interconnected by a high-voltage transmission grid, which can absorb a minimum phase unbalance without difficulty.<sup>2</sup> The northeastern sector of the United States, together with areas in the southwest and northwest, has such a system of power generator and transmission. This sector covers the territory roughly north of the southern limit of the Tennessee Valley Authority area and the territory east of the Mississippi River. It practically coincides with the ICC Eastern District of Class I railways.

### COMMUNICATION

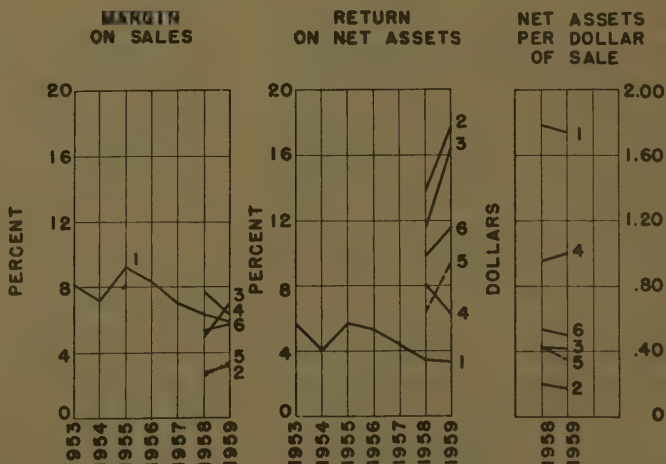
Automation can be used effectively to reduce operating cost. As applied to railways, in addition to the modern push-

button classification yards, this involves centralized data processing and electronic computer control of such things as payrolls, waybills, invoices, inventory, reservations, accounting data, train make-up, and movements of trains and cars. Transmission of large quantities of data in digital form is required from remote points to the central computer center, with reversed transmission of instructions, again in digital form, from the center to remote points. Envisioned for the future are systems of automatic train operation under control of a central computer center, or central brain, requiring continuous flow of feedback information from trains to the central brain and a reverse flow of instructions from the central brain to trains.

Moreover, the concept of "real-time" operation of railroads presupposes the

Fig. 3. Economic results of transportation company operations

- 1—Class I railroads
- 2—Motor carriers
- 3—Interstate bus companies
- 4—Inland water carriers
- 5—Air transport
- 6—Total manufacturing (for comparison)





ability to determine in advance the various movements of trains and cars required to accommodate the ebb and flow of traffic to be delivered from or to a given railroad. The production of this predetermined information requires complete use of data processing and computer facilities.

Practical automation of railway operating functions calls for a communication system that will provide many digital transmission channels and a limited number of voice channels, at low cost per channel. A microwave system for through-communication would meet this requirement. One that would be suitable for many railroads provides 30 voice channels, each of which can be subdivided into 15 digital channels. This gives a possible maximum of 450 digital channels, or say 10 voice channels and 300 digital channels. Initial capital investment and maintenance and operating costs are low as compared with any other system. Microwave is not affected by inductive interference from the traction electric power supply system. Relay stations could be located about 30 miles apart, which roughly coincides with the desired spacing of power supply points for the electrification system. For through-transmission of information, microwave appears to be ideal as a part of the technological package.

There remains, however, the problem of a local communication system which will make possible end-to-end communication on trains, between trains, between trains and wayside stations, and through communication between a control center and wayside stations, and between the control center and trains. Very-high-frequency (vhf) radiobroadcasting from microwave relay stations could be used. A more desirable system, however, would be the use of vhf signals, carried by the conductors of the electrification power supply system, and inductively coupled to moving trains and wayside locations. There are two known systems of this type, both of which probably would require further development for the proposed application. Any such system must provide numerous digital channels and a few voice channels. It should also make possible the bridging of necessary emergency through communication channels if there is a loss of one microwave repeater station.

A high-capacity communication system, having a low initial capital cost and free of inductive interference from the electrification power system, can be provided, and it should produce maintenance and operating savings over existing sys-

tems sufficient to justify its initial cost. A microwave system, coupled with a vhf local system, appears proper for inclusion in the technological package.

#### CENTRALIZED TRAFFIC CONTROL

Taken as a group, railroads in the northeast sector of the United States have more track miles than needed to handle their maximum traffic. Consequently they are burdened with excessive cost in the maintenance of way accounts. Centralized traffic control (CTC) is a well-known system by which train movement can be controlled from a single location. A dispatcher at that point can operate track switches and wayside signals up to distances as remote as 600 miles from his office. This system will allow elimination of trackage and of manually controlled interlockers, thereby producing savings on maintenance and operating costs.

In general, a single-track railroad with CTC has 80% the capacity of a 2-track railroad without it. In many cases, a 2-track road can be reduced to a single-track by installing CTC, and it will still have adequate capacity for handling existing traffic. As will be seen later, electrification will increase the speed of train movement, which will permit reduction in trackage without decreasing capacity. Similarly, 4-track and 3-track railroads can have less track mileage. Reduction in the number of tracks also makes possible the establishment of a work roadway for off-track maintenance. The resulting savings will justify the capital expenditure for the installation of CTC. Furthermore, the reduction also cuts the cost of electrification. CTC requires new signaling of the railroad, but its cost is justified by the savings produced. In this way, the signal changes which ordinarily would result from electrification are not burdensome.

CTC qualifies for inclusion in the package because it reduces trackage and produces permanent maintenance and operating savings. When it is combined with higher train speeds resulting from electrification, equivalent or increased track capacity is provided. Off-track maintenance, made possible by eliminating excessive tracks, is another advantage.

#### CONVERTED ELECTRIC LOCOMOTIVES

It is generally considered, at the average annual usage made of road diesel-electric locomotives by the Class I railroads of the United States, that their economic life is 12 to 15 years. At the end of this period the prime mover and its auxiliaries must be replaced. How-

ever, the mechanical parts of the locomotive, the traction motors, and the air-brake system could be used economically for at least another 15 years as part of a diesel locomotive. Basically, the high rate of increase in maintenance cost as the locomotive ages is responsible for its short economic life. An electric locomotive has a rate of increase one-fourth or less than that of a diesel-electric locomotive.<sup>3,4</sup> Therefore, conversion to a rectifier-type electric locomotive is practical, using the still-serviceable traction motors, the mechanical parts, and the air-brake system of the diesel.

By converting these locomotives, in the general range of heavy freight train operation, gross ton-miles per train hour can be nearly doubled, provided speed is not restricted by track curvature or other interference. This is because the electric locomotive is not limited by constant horsepower, as is the diesel, and within the capacity of the traction motors, can exert its maximum and continuous tractive efforts up to much higher speeds. The tractive effort at the maximum locomotive speeds will not differ from those of the locomotive as a diesel-electric unit. Due to the absence of engine vibration and oil contamination, the mechanical parts will last longer, and conversion will cost less than a new diesel-electric unit.

Conversion of diesels to straight electric locomotives qualifies as part of the technological package because it will result in enormous maintenance savings, provide longer engine life, and, in most cases, will add marginal savings through increased train speeds and possibly lower energy cost. It will conserve the use of parts of old diesels and thereby reduce the cost of locomotive replacement to a minimum. Also, more attractive service will ensue.

#### ELECTRIFICATION SYSTEM

The system to be included in the technological package is the single-phase commercial-frequency system, using 25,000 volts on the distribution system of the tracks. Single-phase power is taken directly from the 3-phase high-voltage network with alternate phases supplied to different sections of the railway in order roughly to balance the phase load. Studies made in this country and the experience gained from electrification in France indicate that the residual power unbalance probably will not exceed the limit of tolerance for larger power systems.<sup>1,2</sup> However, individual studies will be necessary in each case. If required, known remedial measures can be applied.



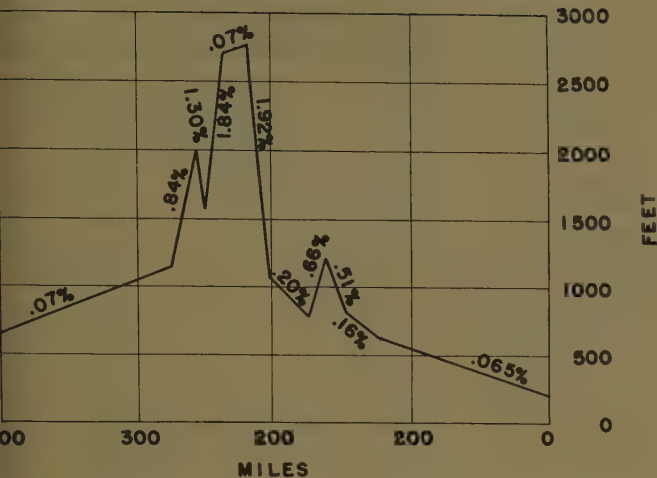


Fig. 4. Assumed average profile

As has been indicated, the central station industry may be willing to sell 25,000-volt power to the railway at the right of way, thus eliminating investment in transmission lines and stepdown substations as part of the electrification system. Comparatively good load factors and diversity factors with existing demands make the railway load attractive for addition to central-station loading. Sixty-cycle commercial frequency makes possible transformers of optimum size and cost as applied in the power stepdown substations and on the rectifier-type locomotives.

The 25,000-volt distribution voltage reduces the equivalent copper and the weight of the catenary system for normal voltage regulation and increases spacing between the power substations. The reduced catenary weight makes possible smaller and lighter poles and foundations, and a consequent over-all reduction in system costs. With power feeds at 25- to 30-mile spacing, it is proposed that master breakers at stepdown stations would control the entire section between feedpoints in order to remove power in the event of a short-circuit fault. Disconnect switches at intermediate sectionalizing stations would localize the fault; then power could be reapplied by closing the master breakers. This could be done in a few seconds without seriously disturbing train operation. For heavily loaded or unusually long sections, auxiliary feeders could be used, or a 3-wire system employing intermediate autotransformers with track center as the mid-point between feeders and valley.

This electrification system qualifies for inclusion in the technological package and that its first cost will be less than any other known system, and its operating and maintenance cost will be relatively low.

### Economic Study

The economic advantages of the proposed technological package have been determined through a comparative economic study of all major costs, both with dieselization and with electrification. The existing cost level was determined for the year preceding conversion, and annual costs then derived for each of the 12 years following. The immediate effects of conversion by installation of the technological package could thus be shown.

#### BASES

The generalized study was based upon data, at a constant 1958 dollar value, from ICC published information and from other sources, applying to a typical Eastern District railroad operating from the eastern seaboard inland. Data were developed for four traffic levels utilizing a given railroad plant layout. The general bases of the study are as follows:

1. The existing railroad plant layout consists of a 2-track railroad having an average profile as shown in Fig. 4. There is some third track on the heavy-grade sections to accommodate movements of pusher locomotives. It was assumed that by use of CTC this basic plant could be converted to a single-track railroad with relatively frequently spaced passing sidings. The average spacing of these sidings varied with traffic level.

2. Operating data were developed on the basis of standard freight trains having the same trailing tonnage and the same number of road locomotives. All Eastern District railroads operating inland from the eastern seaboard have a

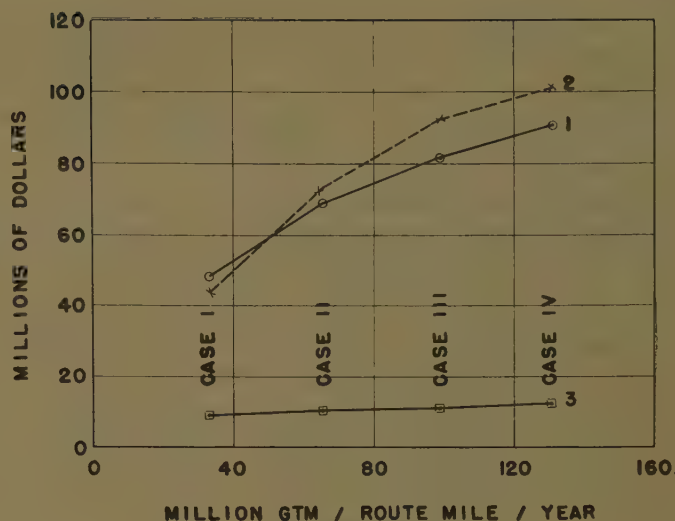


Fig. 5. Estimated net investment

- 1—Package with electrification
- 2—Package with dieselization
- 3—Communication and CTC, signal portion of package

considerably higher volume of revenue tonnage eastbound than westbound. This characteristic was reflected in the operating data in that fewer trains with lower trailing tonnage, but with the number of locomotive units required to balance the movement of motive power, made up the westbound train movement. Gross tons per train and the number of road and pusher locomotives were determined on the basis of ruling grades instead of the average grades which are shown in Fig. 4.

3. For simplification, it was assumed that a large fleet of diesel-electric locomotives, of 1,250 rhp (rail horsepower) rating, all 19 years of age, was available for electrification. The existing locomotives would be operating in their last depreciable year and at maximum maintenance cost. This cost was not based on system fleet averages; the average cost at the 16th year of life was used, as shown by published studies.<sup>3</sup> It was also assumed that the mechanical parts, traction motors, and air-brake equipment of the converted locomotive would have a life of at least 15 years, at which time the packaged electric equipment could be used to convert another fully depreciated diesel-electric locomotive.

4. Another assumption was that 30% of the capital cost for installation of a new communication system and the required CTC system would be met from depreciation funds on hand and the remaining cost from funds derived from the sale of bonds on the open market. It was further assumed that existing diesel-



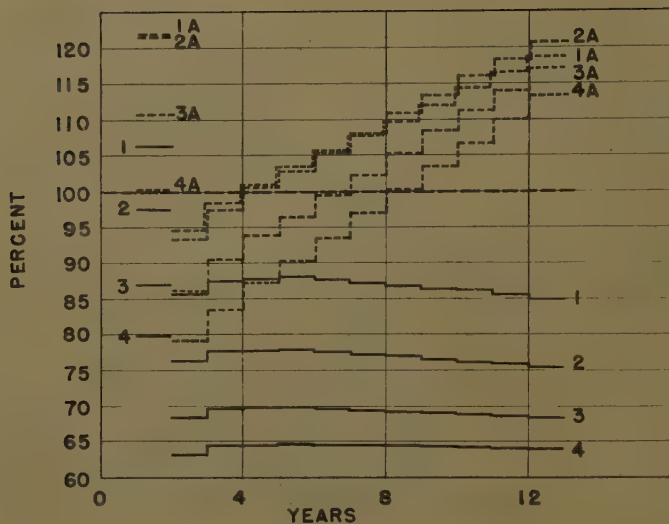


Fig. 6. Comparison of estimated existing cost level (100%) with estimated annual costs, including payments on capital, after conversion

1, 2, 3, 4—Cases I, II, III, and IV, respectively, with electrification  
1A, 2A, 3A, 4A—Cases I, II, III, and IV, respectively, with dieselization

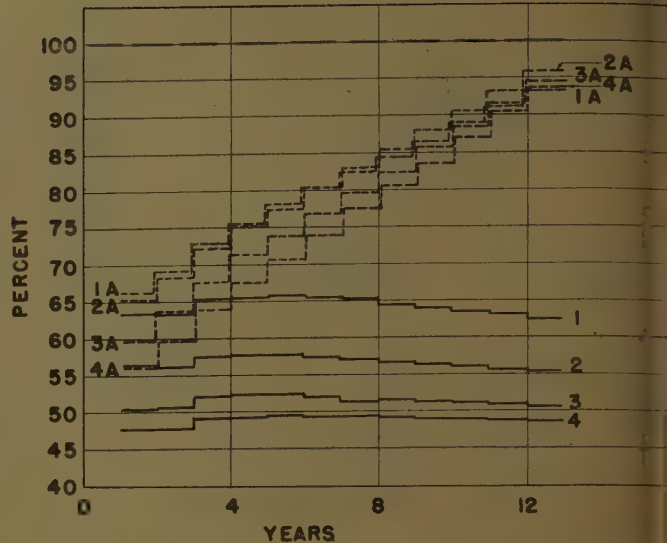


Fig. 7. Comparison of estimated existing cost level (100%) with estimated annual costs, not including payments on capital, after conversion

1, 2, 3, 4—Cases I, II, III, and IV, respectively, with electrification  
1A, 2A, 3A, 4A—Cases I, II, III, and IV, respectively, with dieselization

electric locomotives would either be replaced by diesel-electric locomotives of a somewhat higher horsepower rating or would be converted to electric locomotives, and that the capital cost for this would be met by the sale of equipment trust certificates without recourse to depreciation funds. The cost of the new substations and distribution circuits required for electrification would be met by sale of bonds. In the case of both bond issues, sinking funds would be set up and the interest earnings of those funds would be credited against the interest cost on the bonds.

5. The existing cost level was established by determining charges for:

(a). Depreciation and property taxes on the existing communication and signal systems, together with depreciation for the last depreciable year on existing locomotives.

(b). Operation and maintenance of the existing communication and signal systems, maintenance of existing trackage, locomotive maintenance, engine house expense, and lubrication of locomotives.

(It will be noted that the existing cost level includes neither interest charges nor payments on capital.)

6. The total successive annual costs were established as the summation of:

(a). Sinking fund payments on bond issues and amortization payments on equipment trust certificates.

(b). Net interest on bond issues and interest on the outstanding balance for equipment trust certificates.

(c). Depreciation charges on the new communication and CTC systems, the electric power system, and new diesel-electric locomotives or the converted locomotives.

(d). Property taxes on the new communication, CTC, and electric power supply systems.

(e). Operating costs for the same items listed in (d).

(f). Maintenance costs for the same items as in (d) and (e) and for the revised track layout, and locomotives—either diesel-electric or converted electric—together with engine house expense and locomotive lubrication.

Unit values used in establishing both

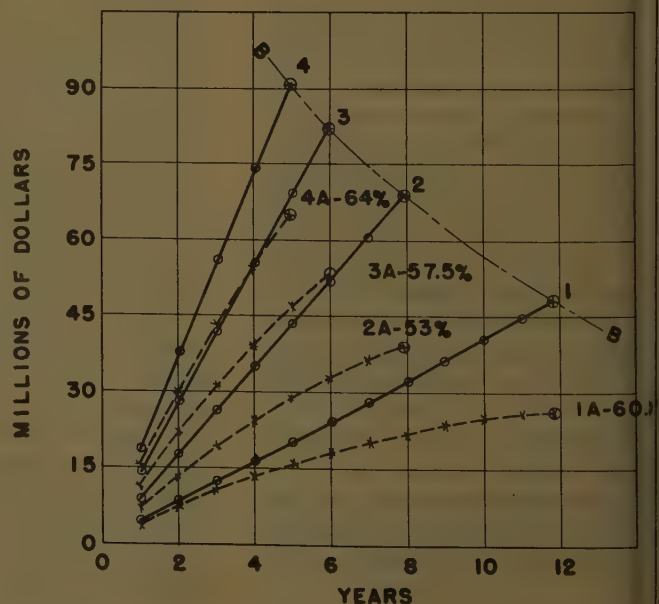
existing cost levels and annual costs after conversion are believed to be conservative. They are available from author upon application.

#### ANALYSIS

Fig. 5 shows the required investment for the four cases studied, including the proportion going for the new communication and CTC systems. These systems require a relatively constant investment whereas the totals show increases which occur at decreasing rates for the heavy traffic cases. In all except case I, total investment for the technology package utilizing railway electrification

Fig. 8. Cumulative savings

1, 2, 3, 4—Cases I, II, III, and IV, respectively, with electrification. 1A, 2A, 3A, 4A—Cases I, II, III, and IV, respectively, with dieselization. B-B—Curve of 100% accumulation of net capital investment for package with electrification; percentage figures indicate portion of net capital investment which would accumulate in time shown for package with dieselization





**Table I. Proportion of Package Represented by New Communication and CTC Systems**

Cases	Type of Loco-motive	Proportion of Capital Investment	Proportion* of 100% Accumulation
I—Diesel	.....	20.0.....	0
	Electric.....	18.3.....	24.2
II—Diesel	.....	14.2.....	0
	Electric.....	15.2.....	15.2
III—Diesel	.....	12.1.....	0
	Electric.....	13.5.....	13.9
IV—Diesel	.....	12.0.....	0
	Electric.....	13.7.....	12.1

\* The dollar proportion in each case is approximately \$11,000,000.

less than that for the package utilizing dieselization. This is attributed to lower cost of conversion as compared with purchasing new diesel-electric locomotives, and to a somewhat smaller fleet of electric locomotives. As the result of higher train speeds and consequent better locomotive utilization, fewer electric locomotives are needed.

One method of analyzing the results obtained by conversion is to compare annual costs, including payments on capital. Since dollar values are different for each case studied, Fig. 6 has been prepared to show the relationship of the four cases in percentages of the existing cost level. This chart clearly demonstrates that, except for the first year when large first payments were made, the technological package with electrification has permanently stabilized annual costs below the existing level.

The package with dieselization is markedly different in that, except for the first year, costs initially drop below but not as far below the existing level as in the case of electrification, and then rapidly rise to well above that level. Costs below the existing level mean increased profits, while those above it signify decreased profits. The striking difference in the two packages arises from comparatively rapid rise in diesel maintenance costs, coupled with greater charges for financing and depreciation which result from higher unit cost per locomotive, shorter economic life, and a larger required fleet.

The economic results of conversion by the technological package can be shown by studying the existing cost level relationship to successive annual costs after conversion, but not including payments of capital; see Fig. 7. The difference between the two costs represents savings which could be effected and which can be expressed cumulatively in terms of dollars for each case. This was done in Fig. 8.

This chart shows the basic divergence in trends in a different manner. In each case, the rate of rise in cumulative savings is greater, and practically uniform, for the package with electrification; but has a drooping characteristic at a lower level for the package with dieselization. In each of the four cases, 100% accumulation of total capital costs for the technological package with electrification is achieved long before such accumulation for the package with dieselization would be attained.

Division of the technological package between new communication and CTC systems and motive-power requirements is indicated in Table I. The communication and CTC portion of the package is relatively small in both capital investment and cumulative savings. The bulk of capital investment and cumulative savings are accounted for by the motive-power portion of the package in each case studied. In this connection, the capital cost of the electrification substation and distribution systems accounts for 40% in case I and 23% in case IV of the total required motive-power investment.

The new diesel-electric locomotives were assumed to have a moderately higher horsepower rating than those previously used. By converting existing diesels to electric locomotives, the continuous rhp would be more than doubled because, at the low speeds at which the continuous tractive effort of the diesel-electric locomotive is exerted, the traction motors are grossly underloaded as to applied voltage. Conversion to electric locomotives would make possible full-capacity use of traction motors at the continuous rating. These increased continuous horsepower ratings mean higher average train speeds, as indicated in Table II.

Electrification will make possible an average round-trip speed, neglecting stops and turn-around time, 25% greater than obtained with existing diesels. The full advantage of the electric locomotive is brought out in the data showing the average speed ascending the 1.92% grade, where use of electric locomotives would increase the speed up the grade approximately 2.5 times more than that now obtained.

## DISCUSSION OF RESULTS

The economic advantage of installing the technological package with electrification of the railroad is clearly indicated for the conditions assumed. Even when payments on capital are included, there is a reduction of cost below the prevailing level and in all cases the level remained approximately constant over the 12-year period.

In the case of the technological package with dieselization, there is also a drop below the existing cost level, but to a lesser degree, followed by a pronounced upward trend in the annual cost. Therefore, when payments on capital are included, the annual cost level soon exceeds the existing one. This is true for all four cases studied.

The study also indicates that cumulative savings for the package with electrification would equal 100% of the net capital investment in 5 to 12 years, depending upon the traffic level. Cumulative savings for such periods materially exceed those of the package with dieselization.

The results of a general study of this nature can be altered by refinement of the estimates as to (1) capital investment cost required, (2) existing cost level to be used as a reference, (3) annual costs which will be incurred after installation of the technological package, and (4) operation over different profiles. Other cases could be studied involving an existing single-track, 3-track, or 4-track railroad, with average profiles and traffic levels differing from those used in this case. Also, definitive studies can be made of individual situations, using factual data rather than estimates.

This analysis shows, however, that in all studies involving relatively heavy traffic the same trends will prevail when comparing dieselization and electrification. The effects of lower capital costs, smaller locomotive fleet, longer economic life, and radically lower rate of rise in maintenance costs for electric locomotives overwhelmingly dictate a lower and stabilized cost level as compared with that of the technological package with dieselization. In all such cases, an economic advantage will be realized by installing the package with railway electrification. The rate of rise of cumulative savings and the time periods required for 100% accumulation of capital requirements will vary, but the basic relationship between the packages will be the same.

**Table II. Train Speeds and Locomotive Horsepower**

Type of Locomotive	Average Round-Trip Speed, %	Speed Ascending 1.92% Grade, %	RHP per Locomotive, %
Existing diesel-electric.....	100.0.....	100.0.....	100.0
New diesel-electric.....	107.5.....	122.0.....	116.0
Converted electric.....	125.0.....	256.0*	212.0

\* At maximum short-time horsepower rating.



Conclusions

This study demonstrates that the installation of the technological package with railway electrification would meet the basic needs of the railroads in the Eastern District, as stated previously, by reduction in operating and maintenance costs with a minimum of capital expenditure.

Discussion

L. W. Birch (Ohio Brass Company, Mansfield, Ohio): In my old psychology book there are two types of escape mechanisms, vividly described, that stand out in my memory. The first escape mechanism deals with the "scapegoat" and the second with projection—accusing the other fellow of your own faults. Why have I mentioned these two escape mechanisms? I have done this because the overhead catenary distribution systems are usually considered extremely costly and have, many times, been offered as a reason why an attempt to electrify a railroad has failed. Let us see if this reasoning is right.

Some readers may be familiar with the work of the Joint Committee on Railroad Electrification. This committee, which consisted of members from railroads, the coal industry, power companies, cable manufacturers, electric equipment manufacturers, and railway consultants, was set up to study the future of railroad electrification in the United States and Canada. The committee, active in the early 1950's, spent almost 2 years in this study. Three contact voltages were involved: (1) 3,000 volts direct current; (2) 12,000 volts alternating current; and (3) 25,000 volts alternating current, mentioned in the paper.

In 1950 and 1951, the use of 25-kv trolley voltage was foreign to most of those participating in the study. However, a few of us had helped design and install the 17-mile electrification of Henry Ford's Detroit, Toledo & Ironton Railroad in 1925-26. The trolley voltage of this system was 22 kv. In 1944-45, several of the Land Transportation Committee members worked on a New York Central estimate, involving electrification of the Central from Harmon to Buffalo, N. Y. This study was based on the same trolley voltages that had been considered in 1950-51 by the Joint Committee on Railroad Electrification.

Let us consider the New York Central costs and the Joint Committee's costs for 25-kv overhead, not including certain standard percentages such as contractor's profits, insurance, and interest. We will compare only the catenary overhead system. If the 1950 cost index for the 25-kv catenary overhead, without transmission lines, is 100, then the 1951 Joint Committee's index for the same voltage

Also, to assist in attracting additional traffic, greatly increased speed for freight trains would be required.

References

1. FRENCH TECHNICAL ADVANCES IN THE FIELD OF RAILROAD ELECTRIFICATION, F. Nouvion. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 79, Sept. 1960, pp. 241-48.

2. TECHNICAL ASPECTS OF PROVIDING SERVICE TO SINGLE-PHASE 60-CYCLE RAILROAD LOADS, T. J. Nagel, A. F. Gabrielle. *Ibid.*, vol. 77, July 1958, pp. 172-76.  
3. LOCOMOTIVE REPAIR COSTS AND THEIR ECONOMIC MEANING TO THE RAILWAYS OF THE UNITED STATES, H. F. Brown. *Ibid.*, pt. I (*Communication and Electronics*), vol. 80, Sept. 1961, pp. 209-16.  
4. VIRGINIAN RAILWAY MOTOR-GENERATOR ELECTRIC LOCOMOTIVE MAINTENANCE COSTS, T. W. Perkinson. *Ibid.*, pt. II (*Applications and Industry*), vol. 79, Mar. 1960, pp. 33-35.

and the same system is 84, and a review of today's costs indicate that the cost index for a 25-kv catenary system is approximately 130—not an excessive increase over the 1945 index.

It is doubtful if the overhead cost increases are as high, percentagewise, as some of the other costs such as equipment or signaling. Also, there have been many improvements in the distribution art since 1945. Based on these indexes for catenary overhead distribution costs, there is little cause to ask, "Are the overhead distribution costs retarding railway electrification?"—the subject of a paper which I presented at the 1948 AIEE Winter General Meeting.

B. A. Ross (American Electric Power Service Corporation, New York, N. Y.): The author's analysis of the economics of a proposed railroad improvement package, including commercial-frequency electrification, is both interesting and timely.

The utility industry is well able to handle the railroad loads with a minimum of new facilities. Studies indicate that existing facilities can, in most areas, tolerate the phase unbalances which might result from the railway loads. While the utility can deliver energy at the catenary voltage as indicated, it should be remembered that, in most areas, only the major subtransmission circuits (i.e., those circuits in the 100- to 150-kv voltage range) could tolerate the phase unbalances of the single-phase railway load without affecting the equipment of other utility customers. This would probably require, in some areas, limited extensions in the utility transmission system, or slight changes in the number and location of the railway substation facilities.

An additional economic factor favoring electrification is the ability of the utility industry to hold, or perhaps even decrease, the cost of electric energy, whereas future oil prices will, in all probability, continue to rise. The electric utility's ability to control its future energy prices is influenced by the fact that it can: (1) utilize the most economic energy source or combination of sources; (2) employ either conventional or unconventional methods of energy conversion with little restriction as to plant location or size; (3) still further increase the efficiency of energy conversion by using larger generating units, new methods of conversion, and improved techniques;

and (4) improve over-all electric system efficiency and reliability by utilizing higher transmission voltages, higher capacity, and more efficient station equipment, as well as improved relaying and operating methods.

H. C. Cross: In making the economic study referred to by Mr. Ross, the element of power cost was omitted, since today's cost of electrical and diesel fuel, delivered at the locomotive, are considered roughly equal. Significant differences may be found in localized geographical areas. Since the electric utility industry is able to control future energy prices, this element should have thorough consideration when a definitive feasibility study is made. Also noteworthy is Mr. Ross' opinion that phase unbalance can be tolerated in most areas, thus assuming the industry's capability of handling "the railroad loads with a minimum of new facilities."

The remarks by Mr. Birch indicate that the estimated cost of a 25-kv catenary system has not increased during the last 15 years in proportion to other elements of electrification cost. In this study, the estimate is approximately the same as that of the CTC system on a track-mile basis. For a project such as assumed, the reduction in track miles by use of CTC eliminates approximately 35% of the catenary system cost required for a 2-track railroad. Perfection of off-track installation methods might provide reductions in unit costs below these estimates.

In addition to cost, the fixed property investment required for a catenary system has deterred adoption of electrification. However, today's railroads are making just such investments by adopting CTC in order to obtain operating and maintenance savings. The over-all savings resulting from use of the technological package adequately justify the catenary system.

The most significant factor involved in the study is maintenance cost of diesel electric locomotives. Following presentation of the paper, the author has received an indication that railway management may not be aware of the extent to which this cost rises with advancing age of their locomotives. Close study of this trend by individual railroads probably will disclose, in many cases, a steeper rate of rise than that estimated in the study. This factor is crucial with respect to economic justification of railway electrification.



# On the Optimum Synthesis of Random Sampling Multipole Filters with Stationary Inputs

H. C. HSIEH  
NONMEMBER AIEE

IN MANY practical engineering problems the sampling intervals for the system are not fixed. They are actually random variables. For example, in radar with nominal scanning period  $T$ , occasionally some samples of the return signals are absent because of noise or other interference. It is referred to as the "missed sample" problem in radar.

This kind of system has been studied to a certain extent by several previous investigators. Kalman considered the optimum synthesis of a random sampling system with step input and quadratic error criterion.<sup>1</sup> His synthesis procedure is based on the technique of dynamic programming. He made detailed studies on the stability of such systems.

Bergen investigated the minimum mean-square synthesis of a randomly sampled system with one stationary stochastic input.<sup>2</sup> The evaluation of power spectral density after sampling in terms of that before sampling and the generating function of the random sampling process has been obtained. His synthesis procedure is based upon the standard Wiener's spectral-density factorization technique.<sup>3</sup> It has been pointed out that, with the assumption of rational functions in  $s$  for the spectral densities before sampling, straightforward factorization of the spectral densities after sampling can be carried out for two common and important classes of sampling.

The purpose of this paper is to obtain the design equations for an optimum random sampling multipole filter with  $n$  stationary inputs and  $m$  outputs. This

optimum system is also optimum in the sense of Wiener. The system is assumed to be linear and time invariant. A set of integral equations which the optimum weighting functions must satisfy is obtained. Transform method is then used to solve these equations.

## Formulation of the Filtering Problem

### DESCRIPTION OF RANDOM SAMPLING PROCESSES AND SYSTEM INPUTS

The random sampling of an input series means that the samples of this input are taken at random time  $t_i$ . The  $i$ th interval  $T_i$  between the sampling times  $t_i$  and  $t_{i+1}$  is defined as:

$$T_i = t_{i+1} - t_i \quad (1)$$

Instead of working directly with the sampling time  $t_i$ , it is convenient to assume that these sampling intervals  $T_i$  are statistically independent random variables with a first-order probability density function  $f_1(T)$ . In other words,

$$Pr(T < T_i < T + dT) = f_1(T) dT \text{ for all } i \quad (2)$$

This density function is either continuous or discrete.

Let the time series  $u(t)$  be a sequence of very narrow pulses with width  $\gamma$  and of unity amplitude, which occur at random times. The intervals  $T_k$  between the leading edges of successive pulses are statistically independent variables with the same first-order probability density function  $f_1(T)$ . This sampling process

is thus stationary. The random pulse modulation of the input may then be described by

$$i_u(t) = \frac{1}{\bar{u}} u(t) i(t) \quad (3A)$$

Here  $\bar{u}$  is the ensemble average of  $u(t)$  and is arbitrarily introduced in the equation in order to simplify the subsequent manipulations. This sequence of random pulse modulation of the input is shown in Fig. 1. The randomly sampled input  $i^*(t)$ , which is usually regarded as the impulse modulation of the continuous input can then be obtained as the pulse width approaches zero. Hence:

$$\begin{aligned} i^*(t) &= \lim_{\gamma \rightarrow 0} \frac{1}{\bar{u}} u(t) i(t) \\ &= \frac{1}{\bar{u}} u^*(t) i(t) \end{aligned} \quad (3B)$$

where the asterisk denotes the sampled quantity and the bar, the statistical average.

The system diagram is shown in Fig. 2. The input before sampling at each terminal of the multipole filter consists of stationary random signal and noise. Thus:

$$i_k(t) = S_k(t) + N_k(t) \quad (4)$$

These processes have zero means and known correlation functions. In dealing with multipole systems, it is reasonable to assume that the random sampling process of each terminal is statistically independent of those of the others and also of the input processes. Hence, for the pulse modulation:

$$i_{uk}(t) = \frac{1}{\bar{u}_k} u_k(t) i_k(t) \quad (5)$$

and

$$\begin{aligned} \overline{u_{k'}(t) u_k(t + \tau)} &= \overline{u_{k'}} \overline{u_k} \text{ for } k' \neq k \\ &= \phi_{u_k u_k}(\tau) \text{ for } k' = k \end{aligned} \quad (6)$$

where  $\phi_{u_k u_k}(\tau)$  is the autocorrelation function of the sampling process. It can easily be shown that, with further assump-

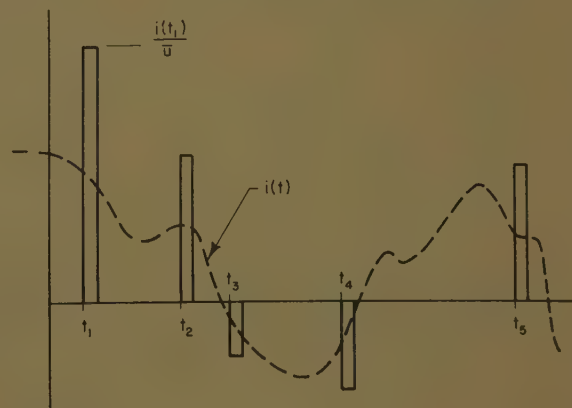


Fig. 1. Pulse sequence  $i_u(t)$  of randomly sampled input

Paper 61-802, recommended by the AIEE Basic Sciences Committee of the Science and Electronics Division and approved by the AIEE Technical Operations Department for presentation at the AIEE Summer General Meeting, Ithaca, N. Y., June 18-23, 1961. Manuscript submitted February 6, 1961; made available for printing April 7, 1961.

H. C. HSIEH is with the University of California, Los Angeles, Calif.

The author expresses his gratitude to Dr. C. T. Leondes of the University of California at Los Angeles, for his many helpful suggestions. This research was supported by the U. S. Air Force under Contract no. AF49(638)-438, monitored by the Air Force Office of Scientific Research of the Air Research and Development Command.



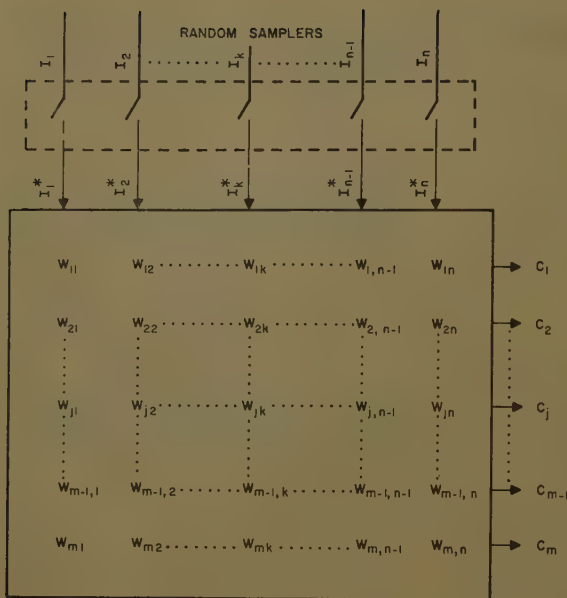
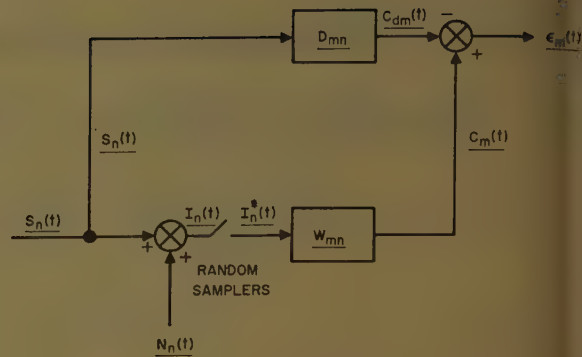


Fig. 2 (left).  
Block diagram for  
 $n+m$  pole filter

Fig. 3 (right).  
Error generation  
diagram



tion of ergodicity for the process, the average value of  $u_k(t)$  is:

$$\bar{u}_k = \frac{\gamma}{T_k} \quad (7)$$

where  $T_k$  is the average sampling period. The derivation of  $\phi_{u_k u_k}(\tau)$  and its limiting case  $\phi_{u_k u_k}^*(\tau)$  is given in reference 2. The latter is expressed by:

$$\phi_{u_k u_k}^*(\tau) = \bar{u}_k \gamma \left[ \delta(\tau) + \sum_{n=1}^{\infty} \int_{T_1}^{T_1+T_n} \dots \int_{T_n}^{T_n+T_n} \delta(\tau - T_1 - \dots - T_n) f_{n_k}(T_1 \dots T_n) dT_1 \dots dT_n \right] \text{ for } \tau \geq 0 \quad (8)$$

where  $f_{n_k}$  is the  $n$ th order probability density function.

#### DERIVATION OF SYSTEM EQUATIONS

In deriving the system equations, pulse modulation is used at the beginning. The system errors are defined, as usual, as the difference between actual outputs and desired outputs; see Fig. 3. Thus:

$$\begin{aligned} e_j(t) &= C_j(t) - C_{dj}(t) \\ &= \sum_{k=1}^n \int_0^\infty W_{jk}(\tau) i_{u_k}(t-\tau) d\tau - \\ &\quad \sum_{k=1}^n \int_{-\infty}^\infty D_{jk}(\tau) S_k(t-\tau) d\tau \\ &= \sum_{k=1}^n \int_0^\infty W_{jk}(\tau) [S_k(t-\tau) + \\ &\quad N_k(t-\tau)] \frac{u_k(t-\tau)}{\bar{u}_k} d\tau - \\ &\quad \sum_{k=1}^n \int_{-\infty}^\infty D_{jk}(\tau) S_k(t-\tau) d\tau \end{aligned} \quad j=1, 2, \dots, m \quad (9)$$

and

$$\begin{aligned} \bar{e}_j^2 &= \sum_{k'=1}^n \sum_{k=1}^n \int_0^\infty \int_0^\infty W_{jk'}(\tau_1) W_{jk}(\tau_2) d\tau_1 d\tau_2 \times \\ &\quad \int_0^\infty W_{jk}(\tau_1) \phi_{i_{k'} i_k}^{u^2}(\tau_2 - \tau_1) d\tau_1 - \\ &\quad 2 \sum_{k'=1}^n \sum_{k=1}^n \int_0^\infty W_{jk'}(\tau_2) d\tau_2 \times \\ &\quad \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{i_{k'} S_k}^u(\tau_2 - \tau_1) d\tau_1 + \\ &\quad \sum_{k'=1}^n \sum_{k=1}^n \int_{-\infty}^\infty D_{jk'}(\tau_2) d\tau_2 \times \\ &\quad \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{S_{k'} S_k}(\tau_2 - \tau_1) d\tau_1 \end{aligned} \quad j=1, 2, \dots, m \quad (10)$$

where

$$\begin{aligned} \phi_{i_{k'} i_k}^{u^2}(\tau) &= \phi_{S_{k'} S_k}^{u^2}(\tau) + \phi_{S_{k'} N_k}^{u^2}(\tau) + \\ &\quad \phi_{N_{k'} S_k}^{u^2}(\tau) + \phi_{N_{k'} N_k}^{u^2}(\tau) \\ \phi_{i_{k'} S_k}^u(\tau) &= \phi_{S_{k'} S_k}^u(\tau) + \phi_{N_{k'} S_k}^u(\tau) \end{aligned}$$

are the various correlation functions. The synthesis criterion is, then, to minimize these system mean-square errors.

Investigation is now made of all the correlation functions after sampling. It has been assumed that each of the sampling processes is statistically independent of the others and also of the input processes.

Thus:

$$\begin{aligned} \phi_{S_{k'} S_k}^{u^2}(\tau) &= S_{k'}(t) \frac{u_{k'}(t)}{\bar{u}_{k'}} S_k(t+\tau) \frac{u_k(t+\tau)}{\bar{u}_k} \\ &= \frac{1}{\bar{u}_{k'}^2} \phi_{u_{k'} u_{k'}}(\tau) \phi_{S_{k'} S_k}(\tau) \text{ for } k=k' \\ &= \phi_{S_{k'} S_k}(\tau) \text{ for } k \neq k' \end{aligned} \quad (11)$$

Similarly

$$\begin{aligned} \phi_{S_{k'} N_k}^{u^2}(\tau) &= \frac{1}{\bar{u}_{k'}^2} \phi_{u_{k'} u_k}(\tau) \phi_{S_{k'} N_k}(\tau) \text{ for } k=k' \\ &= \phi_{S_{k'} N_k}(\tau) \text{ for } k \neq k' \end{aligned} \quad (12)$$

On the other hand,

$$\begin{aligned} \phi_{S_{k'} S_k}^u &= \frac{S_{k'}(t) u_{k'}(t) S_k(t+\tau)}{\bar{u}_{k'}} \\ &= \phi_{S_{k'} S_k} \text{ for all } k \text{ and } k' \end{aligned} \quad (13)$$

and

$$\phi_{N_{k'} S_k}^u(\tau) = \phi_{N_{k'} S_k}(\tau) \text{ for all } k \text{ and } k' \quad (14)$$

Hence the following equations are obtained:

$$\phi_{i_{k'} S_k}^u(\tau) = \phi_{i_{k'} S_k}(\tau) \text{ for all } k \text{ and } k' \quad (15)$$

and

$$\begin{aligned} \phi_{i_{k'} i_k}^{u^2}(\tau) &= \phi_{i_{k'} i_k}(\tau) \text{ for all } k \neq k' \\ &= \frac{1}{\bar{u}_{k'}^2} \phi_{u_{k'} u_k}(\tau) \phi_{i_{k'} i_k}(\tau) \text{ for } k=k' \end{aligned} \quad (16)$$

In other words, the cross-correlation functions of the sampled inputs and the cross-correlation functions of the sampled and continuous inputs, are unaffected by the sampling processes. Since these correlation functions do not depend on the nature of the random-pulse trains, it is evident that, as the pulse width  $\tau$  approaches zero, the following results are true:

$$\phi_{i_{k'} S_k}^*(\tau) = \lim_{\gamma \rightarrow 0} \phi_{i_{k'} S_k}^u(\tau) = \phi_{i_{k'} S_k}(\tau) \text{ for all } k \text{ and } k' \quad (17)$$

and

$$\phi_{i_{k'} i_k}^*(\tau) = \lim_{\gamma \rightarrow 0} \phi_{i_{k'} i_k}^{u^2}(\tau) = \phi_{i_{k'} i_k}(\tau) \text{ for all } k \neq k' \quad (18)$$

The autocorrelation functions of the sampled inputs are, on the other hand, unaffected by the sampling processes. Thus as the pulse width approaches zero we have:

$$\begin{aligned} \phi_{i_k i_k}^*(\tau) &= \lim_{\gamma \rightarrow 0} \phi_{i_k i_k}^{u^2}(\tau) \\ &= \frac{1}{\bar{u}_k^2} \phi_{u_k u_k}^*(\tau) \phi_{i_k i_k}(\tau) \end{aligned} \quad (19)$$

The mean-square errors are, in the limiting case:



$$\begin{aligned}
& \sum_{k'=1}^n \sum_{k=1}^n \int_0^\infty W_{jk'}(\tau_2) d\tau_2 \times \\
& \int_0^\infty W_{jk}(\tau_1) \phi_{i_k' i_k}^*(\tau_2 - \tau_1) d\tau_1 - \\
& 2 \sum_{k'=1}^n \sum_{k=1}^n \int_0^\infty W_{jk'}(\tau_2) d\tau_2 \times \\
& \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{i_k' s_k}(\tau_2 - \tau_1) d\tau_1 + \\
& \sum_{k'=1}^n \sum_{k=1}^n \int_{-\infty}^\infty D_{jk'}(\tau_2) d\tau_2 \times \\
& \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{s_k' s_k}(\tau_2 - \tau_1) d\tau_1 \\
& j = 1, 2, \dots, m \quad (20)
\end{aligned}$$

The necessary and sufficient condition for obtaining minimum errors is that the system weighting functions must satisfy the following integral equations:

$$\begin{aligned}
& \sum_{k=1}^n \int_0^\infty W_{jk}(\tau_1) \phi_{i_k' i_k}^*(\tau_2 - \tau_1) d\tau_1 \\
& = \sum_{k=1}^n \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{i_k' s_k}(\tau_2 - \tau_1) d\tau_1 \\
& \text{for } \tau_2 \geq 0 \\
& k' = 1, 2, \dots, n \\
& j = 1, 2, \dots, m \quad (21)
\end{aligned}$$

These are the generalized Wiener-Hopf equations for a random sampling system. The procedure used in deriving these equations may be found elsewhere.<sup>4,5</sup> The minimum mean-square errors are:

$$\begin{aligned}
& (e^j)^2_{\min} = \sum_{k'=1}^n \sum_{k=1}^n \int_{-\infty}^\infty D_{jk'}(\tau_2) d\tau_2 \times \\
& \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{s_k' s_k}(\tau_2 - \tau_1) d\tau_1 - \\
& \sum_{k'=1}^n \sum_{k=1}^n \int_0^\infty W_{jk'}(\tau_2) d\tau_2 \times \\
& \int_0^\infty W_{jk}(\tau_1) \phi_{i_k' i_k}^*(\tau_2 - \tau_1) d\tau_1 \\
& = \sum_{k'=1}^n \sum_{k=1}^n \int_{-\infty}^\infty D_{jk'}(\tau_2) d\tau_2 \times \\
& \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{s_k' s_k}(\tau_2 - \tau_1) d\tau_1 - \\
& \sum_{k'=1}^n \sum_{k=1}^n \int_0^\infty W_{jk'}(\tau_2) d\tau_2 \times \\
& \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{i_k' s_k}(\tau_2 - \tau_1) d\tau_1 \\
& j = 1, 2, \dots, m \quad (22)
\end{aligned}$$

## Minimization Equations for Optimum Filters

### POWER SPECTRAL DENSITIES AFTER RANDOM SAMPLING

The integral equations for the optimal weighting functions can be solved as

usual by employing the transform method. Thus the transformation of the correlations should first be discussed. As was shown in the previous section, the cross-correlation functions are not affected by the sampling process. Therefore, the cross spectral densities, which are the transforms of the cross-correlation functions, are the same before and after the sampling.

The autocorrelation function of the input after sampling is proportional to the product of the autocorrelation function of that input before sampling and the autocorrelation function of the sampling process, as shown in equation 19. Since the autocorrelation function is an even function, it is possible to express it as:

$$\phi_{kk}^*(\tau) = \phi_{k_1 k_1}^*(-\tau) + \phi_{k_2 k_2}^*(\tau) \quad (23)$$

where

$$\phi_{k_1 k_1}^*(\tau) = \phi_{kk}^*(\tau) \text{ for } \tau > 0$$

and

$$\phi_{k_2 k_2}^*(0) = 1/2 \phi_{kk}^*(0)$$

The two-sided Laplace transform of  $\phi_{kk}^*(\tau)$  is then:

$$G_{kk}^*(s) = G_{k_1 k_1}^*(s) + G_{k_2 k_2}^*(-s) \quad (24)$$

where

$$G_{k_1 k_1}^*(s) = L_1[\phi_{k_1 k_1}^*(\tau)]$$

where  $L_1$  denotes the conventional one-sided Laplace transform.

It has been shown by Bergen that  $G_{k_1 k_1}^*(s)$  can be expressed by a complex convolution integral as:

$$G_{k_1 k_1}^*(s) = \frac{\bar{T}_k}{2\pi j} \int_{c-j\infty}^{c+j\infty} \frac{G_{k_1 k_1}(\lambda)}{1 - M_k(s-\lambda)} d\lambda \quad (25)$$

where  $G_{k_1 k_1}(s)$  is the one-sided Laplace transform of the autocorrelation function before sampling and

$$\begin{aligned}
M_k(s) &= \text{generating function of } \{T\} \\
&\text{in complex variable } -s \\
&= E[e^{-sT}] = L_1[f_k(T)] \quad (26)
\end{aligned}$$

The function  $G_{k_1 k_1}(\lambda)$  will have singularities in the left half of the  $\lambda$  plane only and  $1/1 - M_k(s-\lambda)$  will have singularities in the right half of the  $\lambda$  plane only. There exists an analytic strip around the imaginary axis which separates the singularities of these two functions. The path of integration of the integral lies in this strip. The function  $G_{k_1 k_1}^*(s)$  is analytic in the right half of the  $s$  plane.

The spectral-density convolution integral of equation 25 can be evaluated by contour integration due to the analyticity of  $G_{k_1 k_1}(\lambda)$  and  $1/1 - M_k(s-\lambda)$  in the different halves of the  $\lambda$  plane.

Thus:

$$G_{k_1 k_1}^*(s) = \frac{\bar{T}_k}{4\pi j} \int_{\Gamma_L} G_{k_1 k_1}(\lambda) \frac{1 + M_k(s-\lambda)}{1 - M_k(s-\lambda)} d\lambda \quad (27)$$

where  $[1 + M_k(s-\lambda)/1 - M_k(s-\lambda)]$  is analytic in the left half plane. The closed contour  $\Gamma_L$  includes the path parallel to the imaginary axis and a semi-infinite circle in the left half plane. Thus the only singularities enclosed by this contour are those of  $G_{k_1 k_1}(\lambda)$ .

Several spectral densities after sampling which have been evaluated by Bergen, are given in the following paragraphs.<sup>2</sup>

### Purely Random Sampling

If the sampling times are purely random, then the sampling intervals will be statistically independent with the exponential first-order probability density function.

$$\begin{aligned}
f_{ik}(T) &= \frac{1}{\bar{T}_k} e^{-T/\bar{T}_k} & T \geq 0 \\
&= 0 & T < 0
\end{aligned}$$

Thus

$$M_k(s) = \frac{1}{s\bar{T}_k + 1}$$

and

$$\begin{aligned}
G_{k_1 k_1}^*(s) &= \frac{\bar{T}_k}{2\pi j} \int_{c-j\infty}^{c+j\infty} G_{k_1 k_1}(\lambda) \times \\
&\frac{1}{1 - \frac{1}{(s-\lambda)\bar{T}_k + 1}} d\lambda = G_{k_1 k_1}(s) + \frac{\bar{T}_k \phi_{k_1 k_1}(0)}{2}
\end{aligned}$$

Hence

$$G_{kk}^*(s) = G_{kk}(s) + \bar{T}_k \phi_{kk}(0) \quad (28)$$

The effect of purely random sampling is thus equivalent to adding white noise to the original spectral density.

### Missed Samples Case

For a system with a nominal sampling interval of  $T_0$ , if it is assumed that the misses occur independently and that the probability of a miss is  $p_k$ , then:

$$f_{ik}(T) = \sum_{n=1}^{\infty} (1-p_k) p_k^{n-1} \delta(T - nT_0)$$

and

$$M_k(s) = \frac{(1-p_k) e^{-sT_0}}{1 - p_k e^{-sT_0}}$$

It follows that:

$$\begin{aligned}
G_{kk}^*(s) &= \frac{T_0}{2\pi j} \int_{\Gamma_L} G_{kk}(\lambda) \times \\
&\frac{1 - e^{2\lambda T_0}}{[1 - e^{-(s-\lambda)T_0}][1 - e^{(s+\lambda)T_0}]} d\lambda + \\
&\frac{p_k T_0}{1 - p_k} \phi_{kk}(0) \quad (29)
\end{aligned}$$



For periodic sampling  $p_k=0$ , then:

$$G_{kk}(s) = \frac{T_0}{2\pi j} \int_{\Gamma_L} G_{kk}(\lambda) \frac{1 - e^{2\lambda T_0}}{[1 - e^{-(s-\lambda)T_0}][1 - e^{(s+\lambda)T_0}]} d\lambda \quad (30)$$

Hence the effect of misses is equivalent to adding white noise to the spectral density of a sampled-data system operating without misses.

### Spectral Densities of a Markov Signal

Assume the spectral density of a signal before sampling is:

$$G_{S_k S_k}(s) = \frac{2\beta\sigma^2}{-s^2 + \beta^2}$$

Then:

$$G_{S_k S_k}^*(s) = \frac{\sigma^2 \bar{T}_k [1 - M_k(s + \beta) M_k(-s + \beta)]}{[1 - M_k(s + \beta)][1 - M_k(-s + \beta)]} \quad (31)$$

If  $\beta \rightarrow \infty$  and  $\sigma^2$  is bounded, then the sampled Markov signal will approach a white spectral density:

$$G_{S_k S_k}^* = \sigma^2 \bar{T}_k \quad (32)$$

It should be noted that the continuous spectral density before sampling for this limiting case tends to zero.

### TRANSFORMATION OF INTEGRAL EQUATIONS

The integral equations which must be solved are shown to be:

$$\sum_{k=1}^n \int_0^\infty W_{jk}(\tau_1) \phi_{i_k' i_k}^*(\tau_2 - \tau_1) d\tau_1 = \sum_{k=1}^n \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{i_k' S_k}(\tau_2 - \tau_1) d\tau_1$$

for  $\tau_2 \geq 0$   
 $k' = 1, 2, \dots, n$

These equations should be modified before transformation can be applied. Let  $f_{jk'}(\tau)$  be defined as:

$$f_{jk'}(\tau) = 0 \quad \text{for } \tau \geq 0 \quad (33)$$

Then the integral equations can be rewritten as:

$$\sum_{k=1}^n \int_0^\infty W_{jk}(\tau_1) \phi_{i_k' i_k}^*(\tau_2 - \tau_1) d\tau_1 = \sum_{k=1}^n \int_{-\infty}^\infty D_{jk}(\tau_1) \phi_{i_k' S_k}(\tau_2 - \tau_1) d\tau_1 + f_{jk'}(\tau_2)$$

for all values of  $\tau_2$  and  
 $k' = 1, 2, \dots, n$  (34)

With a stable system and the property of the correlation function  $\phi(\tau)$  or  $\phi^*(\tau)$  for large  $\tau$  (i.e.,  $\phi(\tau) \rightarrow 0$  as  $|\tau| \rightarrow \infty$ ), it is quite obvious that  $f_{jk'}(\tau)$  will also approach zero for very large negative  $\tau$ . In general,

$f_{jk'}(\tau)$  will not have any delta functions.

The application of the two-sided Laplace transform in solving equation 34 depends on the existence of an analytic strip common to all the functions being transformed. For a physically realizable and stable system, the transform of the weighting functions, which are defined as transfer functions of the system, can only have singularities in the left half plane. The transform of  $\phi_{i_k' i_k}^*(\tau)$  or  $\phi_{i_k' S_k}(\tau)$  will have an analytic strip around the imaginary axis. Its singularities in the left half plane arise from the part of the correlation function with positive argument and those in the right half plane, from the part of the correlation function with negative argument. The transform of the arbitrarily defined function  $f_{jk'}(\tau)$  can evidently possess singularities in the right half of the plane only. Thus with an appropriate specification of the desired transfer functions, the existence of a common analytic strip is always obtainable. Thus the inverse transform is unique. It is noticed that this analytic strip contains the imaginary axis which can be chosen as the path for inverse transformation. This line integration can always be replaced by contour integration. Closing the contour in the left half plane will give the function for positive  $t$ .

Let the two-sided Laplace transform be taken on both sides of equation 34.

$$\sum_{k=1}^n G_{i_k' i_k}^*(s) Y_{jk}(s) = \sum_{k=1}^n (Y_{jk}(s) G_{i_k' S_k}(s) + F_{jk'}(s))$$

$k' = 1, 2, \dots, n$  (35)

where

$$Y_{jk}(s) = \int_0^\infty W_{jk}(t) e^{-st} dt$$

(System transfer functions) (36)

$$(Y_{jk}(s) = \int_{-\infty}^\infty D_{jk}(t) e^{-st} dt$$

(Ideal transfer functions) (37)

$$G_{i_k' i_k}^*(s) = \int_{-\infty}^\infty \phi_{i_k' i_k}^*(t) e^{-st} dt = G_{i_k' i_k} \text{ for } k \neq k'$$

(Sampled power spectral densities of inputs) (38)

$$G_{i_k' S_k}(s) = \int_{-\infty}^\infty \phi_{i_k' S_k}(t) e^{-st} dt$$

(Continuous cross-power spectral densities of inputs and signals) (39)

$$F_{jk'}(s) = \int_{-\infty}^0 f_{jk'}(t) e^{-st} dt \quad (40)$$

Define

$$N_{jk'}(s) = \sum_{k=1}^n (Y_{jk}(s) G_{i_k' S_k}(s))$$

Then equation 35 can be written as:

$$\sum_{k=1}^n G_{i_k' i_k}^*(s) Y_{jk}(s) = N_{jk'}(s) + F_{jk'}(s)$$

$k' = 1, 2, \dots, n$

In matrix form:

$$\mathbf{G} \mathbf{Y}_j = \mathbf{N}_j + \mathbf{F}_j^- \text{ for } j = 1, 2, \dots, m$$

where the square matrix  $\mathbf{G}$  is identical for all output terminals and is not singular, and the vector  $\mathbf{F}_j^-$  which is unknown at this moment can only have singularities in the right half plane. When  $s = j\omega$ ,  $\mathbf{G}$  is Hermitian. Inverting equation 43, transfer function vectors are:

$$\mathbf{Y}_j = \mathbf{G}^{-1} (\mathbf{N}_j + \mathbf{F}_j^-) = \frac{\mathbf{A}}{|\mathbf{G}|} (\mathbf{N}_j + \mathbf{F}_j^-)$$

In this equation,  $|\mathbf{G}|$  is the determinant of  $\mathbf{G}$  and should not have any root on the imaginary axis.  $\mathbf{A}$  is the adjoint of  $\mathbf{G}$ .

In complex domain, the minimum mean square errors are

$$\begin{aligned} (\epsilon_j^2)_{\min} &= \sum_{k'=1}^n \sum_{k=1}^n \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} \times \\ &\quad [(Y_{jk'})^* - (Y_{jk})^* G_{i_k' S_k} - Y_{jk'} - Y_{jk} G_{i_k' i_k}^*] ds \\ &= \sum_{k'=1}^n \sum_{k=1}^n \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} \times \\ &\quad [(Y_{jk'})^* - (Y_{jk})^* G_{i_k' S_k} - Y_{jk'} - Y_{jk} G_{i_k' i_k}^*] ds \\ &\quad j = 1, 2, \dots, m \end{aligned}$$

Here the integration path lies in the analytic strip common to all these transformed functions and  $Y^-(s) \triangleq Y(-s)$ .

### Solution of Transfer Functions the Optimum Filters

#### FACTORIZATION OF SPECTRAL DENSITY DETERMINANT

The feasibility of using transform methods to solve the integral equation depends upon the factorization of spectral-density determinant  $|\mathbf{G}|$  into two functions  $G^+(s)G^-(s)$ , where  $G^+(s)$  has poles and zeros in the left half plane and  $G^-(s)$  has poles and zeros in the right half plane. Theoretically, this can always be done if the Paley-Weiner criterion is satisfied. Practically, however, the general procedure specified by Wiener is very difficult computationally unless the function is of certain convenient forms.

It is usually assumed that the spectral densities of continuous inputs can be



ed as rational functions in  $s$ . Thus off-diagonal elements of  $|G|$  are all rational functions in  $s$ . However, the elements along the diagonal are the pole-power spectral densities. They may not be rational functions in  $s$ . This adds additional difficulty in effecting the factorization, which is not encountered in simple two-pole systems.

Two important classes of sampling will be considered. In both cases the spectral densities before sampling are assumed to be rational functions in  $s$ . In the first class, the sampling density generating function is rational in  $s$ . This will include, for example, those sampling intervals governed by the gamma density function.<sup>2</sup> Thus the spectral-density determinant, in this case, will be an even function of  $s$  and the factorization can easily be obtained.

In the second class, the sampling interval-density is discrete. This will consist of the case of missed samples and the case where a finite number of possible sampling intervals exists. In either case, the spectral-density after sampling is a function of  $e^{sT_n}$ , thus the determinant  $|G|$  will be a function in both  $s$  and  $e^{sT_n}$ . There is no way to carry out the factorization of this function. In order to overcome this difficulty, the approximation of  $e^{sT_n}$  by a rational function should be used. The Padé approximants (see the Appendix) are particularly useful in this case. The higher the degrees of numerator and denominator polynomials of the approximation, the better the approximation. However, the complexity of the system designed will be increased accordingly. This approach is justified by the fact that signals usually have power concentrated in low frequencies and the synthesized system is low pass. It should be noted that it is usually desirable to choose polynomials of the same degree for the numerator and denominator of the approximant in order to obtain the simplest form for the approximated sampled spectral density.

Thus it has been shown that for these two classes of sampling the power spectral-density determinants will be even functions in  $s$ . The synthesis procedures can then be similar to those for continuous inputs.<sup>4</sup>

#### OPTIMUM FILTERS FOR TWO IMPORTANT CLASSES OF SAMPLING

In the previous section has demonstrated the fact that, for two classes of sampling, sampled spectral densities are rational functions in  $s$  or can be approximated by rational functions in  $s$ . Thus the spectral-density determinant can be easily factored and imposed such that:

$$|G| = G^+(s)G^-(s)$$

where all the poles and zeros of  $G^+(s)$  must lie in the left half plane and  $G^-(s)$ , which is equal to  $G^+(-s)$ , can only have poles and zeros in the right half plane. There should not be any finite zero along the imaginary axis.

The standard procedure in finding the physically realizable transfer functions for this Wiener type of problem can now be adopted. Thus:

$$\mathbf{Y}_j = \frac{\mathbf{A}}{G^+G^-} [\mathbf{N}_j^+ \mathbf{F}_j^-]$$

or

$$G^+ \mathbf{Y}_j = \frac{1}{G^-} \mathbf{A} \mathbf{N}_j + \frac{1}{G^-} \mathbf{A} \mathbf{F}_j^- \quad (46)$$

The matrix on the left side of the equation can have poles only in the left half plane. Hence it follows that

$$G^+ \mathbf{Y}_j = \left\{ \frac{1}{G^-} \mathbf{A} \mathbf{N}_j \right\}^+ + \left\{ \frac{1}{G^-} \mathbf{A} \mathbf{F}_j^- \right\}^+ \quad (47)$$

where  $\{ \}^+$  denotes that part of the function which has poles in the left half plane only. The elements in  $(1/G^-) \mathbf{A} \mathbf{N}_j$  are all known. Therefore:

$$\left\{ \frac{1}{G^-} \mathbf{A} \mathbf{N}_j \right\}^+ = \int_0^\infty \mathbf{g}(t) e^{-st} dt \quad (48)$$

where

$$\mathbf{g}(t) = \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} \frac{1}{G^-} \mathbf{A} \mathbf{N}_j e^{ts} ds$$

The left half plane poles in  $(1/G^-) \mathbf{A} \mathbf{F}_j^-$  can only come from the matrix  $\mathbf{A}$ . Since  $\mathbf{F}_j^-$  is still unknown, the residues or coefficients associated with these poles can not be determined yet. Thus

$$\left\{ \frac{1}{G^-} \mathbf{A} \mathbf{F}_j^- \right\}^+$$

will be a column vector whose  $k$ th element contains poles which are the left half plane poles from the  $k$ th row elements of matrix  $\mathbf{A}$  and can be expressed:

$$\sum_{i, m_i} \frac{C_{im_i}}{(s - \gamma_i)^{m_i}}$$

where  $\gamma_i$  are the left half plane poles, which can be multiple, and  $C_{im_i}$  are the unknown coefficients. A typical transfer function can now be expressed by

$$Y_{jk} = \frac{1}{G^+} \left[ \left\{ \frac{1}{G^-} \sum_{k'=1}^n A_{k'k} N_{jk'} \right\}^+ + \sum_{i, m_i} \frac{C_{im_i}}{(s - \gamma_i)^{m_i}} \right] \quad (49)$$

In order to determine these coefficients,  $Y_{jk}$  are substituted into the original system equations.

It must be true then:

$$\{\mathbf{G} \mathbf{Y}_j\}^+ = \{\mathbf{N}_j\}^+ \quad (50)$$

or, more explicitly:

$$\left\{ \sum_{k=1}^n G_{ik'} Y_{jk}(s) \right\}^+ = \left\{ \sum_{k=1}^n (Y_d)_{jk}(s) G_{ik'}(s) \right\}^+ \quad \text{for } k' = 1, 2, \dots, m$$

Hence a set of independent linear algebraic equations in terms of  $C_{im_i}$  will be obtained. These coefficients can thus be uniquely determined.

In working an actual problem, it is more desirable to modify the expression for  $Y_{jk}$  as given by equation 50 in order to minimize the accumulated computational errors. Since the left half plane poles of the cofactors  $A_{kk'}$  will eventually be cancelled by the poles of  $G^+$ , the transfer functions will take the form of equation 51:

$$Y_{jk}(s) = \frac{\left[ \prod_u (s - a_u) \right] \left[ C_n s^n + C_{n-1} s^{n-1} + \dots + C_0 \right]}{\left[ \prod_v (s - b_v) \right] \left[ \prod_l (s - d_l) \right]} \quad (51)$$

where

$b_v$  = the zeros of  $G^+$

$d_l$  = the left half plane poles of the ideal transfer functions vector  $(\mathbf{Y}_d)_j$  if they are different from those of the spectral densities or have not been cancelled by the poles of  $G^+$

$a_u$  = the poles of  $G^+$ , which have not been cancelled by the poles of the  $k$ th row element of

$$\left\{ \frac{1}{G^-} \mathbf{A} (\mathbf{N}_j^+ \mathbf{F}_j^-) \right\}^+$$

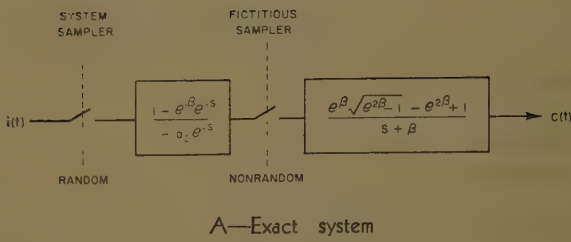
vector, and where  $\{C_n\}$  are the new unknown coefficients which will be a linear combination of the old coefficients  $\{C_{im_i}\}$ . This set of new coefficients is always greater in number than the old ones and, therefore, more simultaneous equations have to be solved. However, with the aid of a digital computer, the solution for these equations can easily be obtained. As a rule, only  $(n-1)$  equations of the  $n$  equations associated with the same pole  $b_v$  of the transfer functions are independent. This procedure will be amplified by referring to example 2, discussed later. The poles of  $Y_{jk}$  are usually simple; however, they can be multiple.

The synthesis procedure for determining the optimum transfer function vector  $\mathbf{Y}_j$  will then be as follows:

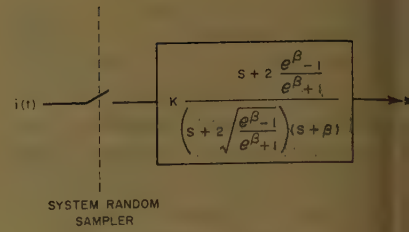


Fig. 4 (left and right). Example 1, probability density function is:

$$\sum_{n=1}^{\infty} 0.8(0.2)^{n-1} \delta(T-n)$$



A—Exact system



B—Approximate system where

$$K = \frac{1+e^{-\beta}}{2} \frac{2 \frac{e^{\beta}-1}{e^{\beta}+1} + \beta}{2 \sqrt{\frac{e^{\beta}-1}{e^{\beta}+1}} + \beta}$$

$$e^{-s} \approx \frac{1 - \frac{1}{2}s}{1 + \frac{1}{2}s}$$

and

$$e^s \approx \frac{1 + \frac{1}{2}s}{1 - \frac{1}{2}s}$$

Then

$$(G_{tt})_1^*(s) = \frac{2e^{\beta}}{e^{\beta}+1} \times \frac{\left(2\sqrt{\frac{e^{\beta}-1}{e^{\beta}+1}} + s\right) \left(2\sqrt{\frac{e^{\beta}-1}{e^{\beta}+1}} - s\right)}{\left(2\frac{e^{\beta}-1}{e^{\beta}+1} + s\right) \left(2\frac{e^{\beta}-1}{e^{\beta}+1} - s\right)}$$

$$Y_1(s) = \frac{1+e^{-\beta}}{2} \times \frac{\left(2\frac{e^{\beta}-1}{e^{\beta}+1} + \beta\right) \left(s + 2\frac{e^{\beta}-1}{e^{\beta}+1}\right)}{\left(2\sqrt{\frac{e^{\beta}-1}{e^{\beta}+1}} + \beta\right) \left(s + 2\sqrt{\frac{e^{\beta}-1}{e^{\beta}+1}}\right)}$$

This approximate system is shown in Fig. 4(B). Due to the approximation used,  $Y_1(s)$  does not satisfy equation exactly.

1. When the spectral densities after sampling are rational functions in  $s$ , the spectral-density determinant  $|G|$  can be factored into  $G^+(s)$  and  $G^-(s)$  directly. When the spectral densities after sampling are functions of  $e^{sT_n}$ , they should be approximated by rational function in  $s$  before factorization.

2. Express the system transfer functions with the form which is given in equation 51.

3. Determine the unknown coefficients by substituting  $Y_j$  into the original system equations.

It is evident that the poles of the transfer functions will consist of two parts. First, all the system transfer functions will have poles which are the zeros of  $G^+(s)$ . These poles, in general, are completely different from the left half plane poles of all the spectral densities. Second, the transfer functions associated with a particular output terminal may contain the left half plane poles from the desired transfer functions vector  $(Y_d)_j$ .

#### ILLUSTRATIVE EXAMPLES

**Example 1.** The validity of using the fractional function approximation of  $e^{sT_n}$  to the optimum synthesis of system will now be shown. A two-pole system will be used so that the mean-square error of the system from the approximate solution can be compared with that from the exact solution.

Let us consider the design of an optimum filter with a nominal sampling period of 1 second and the probability of a miss of 0.2. Assume the signal is Markov with the form

$$G_{SS}(s) = \frac{2\beta}{-s^2 + \beta^2}$$

The noise is white but with a finite mean-square value of 0.6. The signal and noise are uncorrelated. The desired operation is pure filtering.

Thus

$$Y_d(s) = 1$$

In this simple problem, the equation to be solved is:

$$G_{tt}^*(s) Y(s) = Y_d(s) G_{SS}(s) + F^-(s) \quad (52)$$

The sampled spectral density can be found by using equations 29 and 32.

$$G_{tt}^*(s) = G_{SS}^*(s) + G_{NN}^*(s)$$

$$= \left[ \frac{1 - e^{-2\beta}}{(1 - e^{-\beta}e^{-s})(1 - e^{-\beta}e^s)} + \frac{0.2}{0.8} \times 1 \right] + \frac{1}{0.8} \times 0.6$$

$$= e^{-\beta} \frac{2e^{\beta} - e^s - e^{-s}}{(1 - e^{-\beta}e^{-s})(1 - e^{-\beta}e^s)} \quad (53)$$

The numerator can be factored as:

$$2e^{\beta} - e^s - e^{-s} = \frac{1}{a_0} (1 - a_0 e^{-s})(1 - a_0 e^s)$$

where

$$a_0 = e^{\beta} - \sqrt{e^{2\beta} - 1}$$

Thus

$$G^+(s) = b \frac{1 - a_0 e^{-s}}{1 - e^{-\beta}e^{-s}}$$

with

$$b = \sqrt{\frac{e^{-\beta}}{e^{\beta} - \sqrt{e^{2\beta} - 1}}}$$

The exact solution for the optimum transfer function can then be found.

$$Y_0(s) = \frac{1}{G^+} \left\{ \frac{Y_d G_{SS}}{G^-} \right\}^+$$

$$= \frac{1 - e^{-\beta}e^{-s}}{1 - a_0 e^{-s}} \frac{e^{\beta} \sqrt{e^{2\beta} - 1} - e^{2\beta} + 1}{s + \beta} \quad (54)$$

The exact system is shown in Fig. 4(A). The minimum mean-square error is:

$$\epsilon_0^2 = 1 - \frac{(e^{\beta} - e^{-\beta})(e^{\beta} - \sqrt{e^{2\beta} - 1})}{2\beta} \quad (55)$$

Let the approximants chosen be as follows:

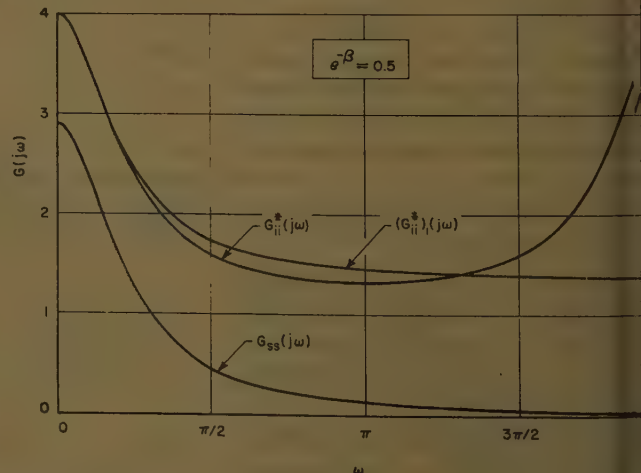


Fig. 5. Power spectral densities in example 1



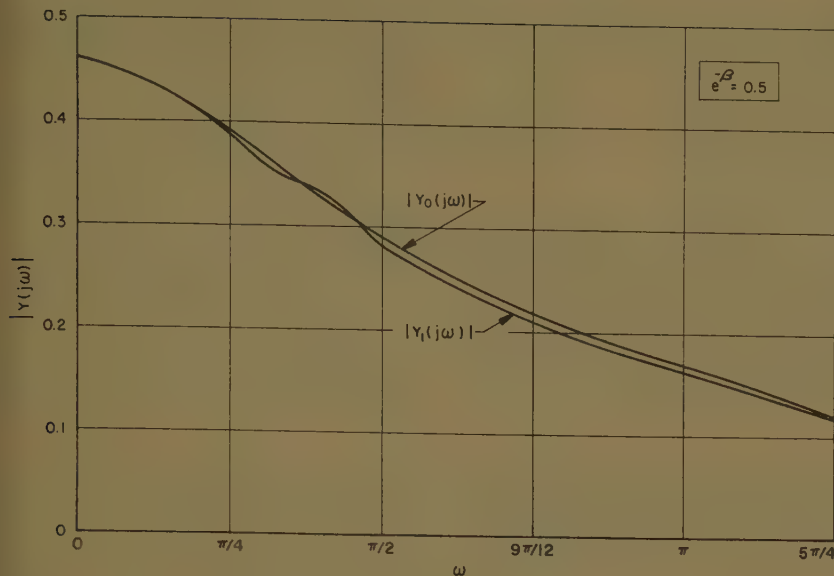


Fig. 6. Gain diagrams in example 1

$$\begin{aligned} \bar{Y}_1^2 = \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} [Y(-s)Y(s)G_{t1}^*(s) - \\ 2Y(-s)Y(s)G_{t2}(s) + \\ Y_d(-s)Y_d(s)G_{ss}(s)] ds \quad (58) \end{aligned}$$

Hence, by using contour integration:

$$\begin{aligned} \bar{Y}_1^2 = \left[ \frac{(1+e^{-\beta})(2a+\beta)}{2(2\sqrt{a}+\beta)} \right]^2 \left\{ \frac{4a^2+2\beta-\beta^2}{2\beta(4a-\beta^2)} - \right. \\ \left. \frac{(4a^2-\beta^2)(10\beta^2-8a)}{8\beta^2(4a-\beta^2)^2} + \right. \\ \left. \frac{(1+e^{-2\beta})(4a^2-\beta^2)}{4\beta(1-e^{-2\beta})(4a-\beta^2)} + \sum_{n=1}^{\infty} \times \right. \\ \left. \frac{n\pi(4a^2-\beta^2+4n^2\pi^2)(4a-\beta^2+4n^2\pi^2)}{16n^2\pi^2-4(a^2-a)} + \right. \\ \left. \frac{2n\pi[(4a-\beta^2+4n^2\pi^2)^2+16n^2\pi^2] \times}{(n^2\pi^2+1)} + \right. \\ \left. \frac{(a^2-a)(2-e^{-(\beta+2\sqrt{a})}-e^{\sqrt{a}-\beta})}{\sqrt{a}(\beta^2-4a)(1-e^{\sqrt{a}-\beta})(1-e^{-(\beta+2\sqrt{a})})} \right\} - \\ \frac{1+e^{-\beta}}{2\beta} \left( \frac{2a+\beta}{2\sqrt{a}+\beta} \right)^2 + 1 \quad (59) \end{aligned}$$

where

$$\Delta = \frac{e^{\beta}-1}{e^{\beta}+1}$$

The infinite series behaves as

$$\sum_{n=1}^{\infty} \frac{1}{2(n^2\pi^2+1)}$$

and, therefore, is convergent.

Fig. 5 shows the spectral densities for real frequencies. It is seen that the approximated sampled spectral density gives a reasonably good fit to the exact one up to frequency  $4\pi/3$ . After that, it deviates from the exact one considerably. However, the magnitude of the signal spectral density  $G_{ss}(j\omega)$  at this frequency range is very small. The filter is certainly of low pass (Fig. 6). Thus,

the error introduced by the approximation is expected to be negligible.

Fig. 7 shows the variation of

$$\bar{\epsilon}_0^2$$

and

$$\bar{\epsilon}_1^2$$

as a function of  $\beta$ . For small  $\beta$ , which implies narrow bandwidth of the spectral densities, these two errors agree very closely. They deviate very slightly as  $\beta$  becomes greater than 1. It can easily be shown that both

$$\bar{\epsilon}_0^2$$

and

$$\bar{\epsilon}_1^2$$

approach zero as  $\beta \rightarrow 0$  and approach 1 as  $\beta \rightarrow \infty$ .

*Example 2.* The synthesis procedure for a multipole filter will now be illus-

trated. For simplicity, a system with two inputs and one output is chosen. The signal portions of the inputs are assumed to be the same. However, they are corrupted by different noises. These noises are not correlated with the signals and are assumed to be white with finite mean-square values. Thus the correlation functions are:

$$\begin{aligned} \phi_{S_1S_1}(\tau) &= \phi_{S_2S_2}(\tau) = e^{-|\tau|} \\ \phi_{N_1N_1}(\tau) &= 0.5 \quad \text{for } \tau=0 \\ &= 0 \quad \text{otherwise} \end{aligned}$$

and

$$\begin{aligned} \phi_{N_2N_2}(\tau) &= 0.2 \quad \text{for } \tau=0 \\ &= 0 \quad \text{otherwise} \end{aligned}$$

The samplers have a nominal sampling period of 1 second but with missing probabilities of 0.1 and 0.2 respectively. The function of the filter is to extract the signal from these two different channels; see Fig. 8.

The power spectral densities for signals before sampling are:

$$G_{S_1S_1}(s) = G_{S_2S_2}(s) = \frac{2}{-s^2+1}$$

and

$$G_{S_1S_2}(s) = G_{S_2S_1}(s) = \frac{2}{-s^2+1}$$

Thus

$$G_{t_1S_1}(s) = G_{t_1S_2}(s) = G_{t_2S_1}(s) = G_{t_2S_2}(s) = \frac{2}{-s^2+1}$$

and

$$G_{t_1t_2}(s) = G_{t_2t_1}(s) = \frac{2}{-s^2+1}$$

The spectral densities for inputs after sampling are:

$$G_{t_1t_1}^*(s) = \frac{1.622-0.245(e^{-s}+e^s)}{1.135-0.368(e^{-s}+e^s)} \quad (60)$$

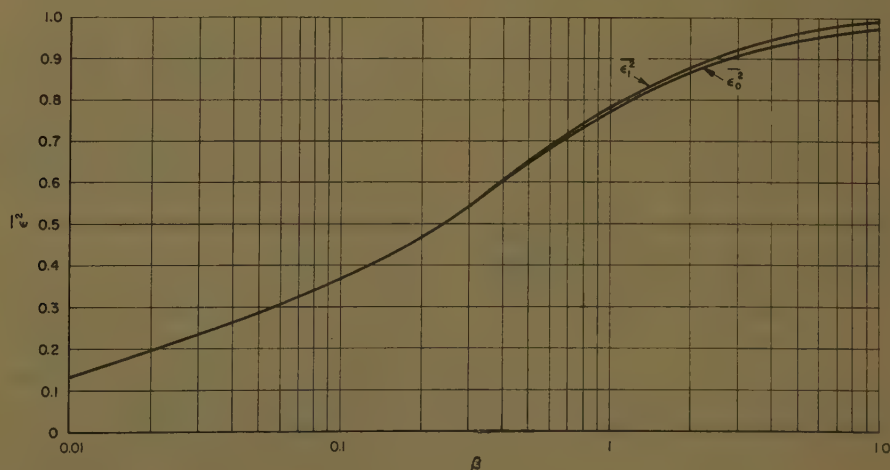


Fig. 7. Minimum mean-square errors in example 1



$$G_{i_2 i_1}^*(s) = \frac{1.433 - 0.184(e^{-s} + e^s)}{1.135 - 0.368(e^{-s} + e^s)} \quad (61)$$

The desired transfer functions are:

$$(Y_d)_{11}(s) = (Y_d)_{12}(s) = \frac{1}{2}$$

The equations to be solved are:

$$\begin{aligned} G_{i_1 i_1}^* Y_{11} + G_{i_1 i_2}^* Y_{12} &= (Y_d)_{11} G_{i_1 s_1} + (Y_d)_{12} G_{i_1 s_2} + F_{11}^- \\ G_{i_2 i_1}^* Y_{11} + G_{i_2 i_2}^* Y_{12} &= (Y_d)_{21} G_{i_2 s_1} + (Y_d)_{22} G_{i_2 s_2} + F_{12}^- \end{aligned} \quad (62)$$

Now let

$$e^{-s} \simeq \frac{1 - \frac{1}{2}s}{1 + \frac{1}{2}s}$$

and

$$e^s \simeq \frac{1 + \frac{1}{2}s}{1 - \frac{1}{2}s}$$

Then

$$G_{i_1 i_1}^*(s) \simeq 1.129 \frac{-s^2 + 2.144}{-s^2 + 0.853} \quad (63)$$

$$G_{i_2 i_2}^*(s) \simeq 0.962 \frac{-s^2 + 2.367}{-s^2 + 0.853} \quad (64)$$

The spectral-density matrix is

$$\mathbf{G} = \begin{bmatrix} 1.129 \frac{-s^2 + 2.144}{-s^2 + 0.853} & 2 \\ 2 & 0.962 \frac{-s^2 + 2.367}{-s^2 + 0.853} \end{bmatrix} \frac{1}{1-s^2} \quad (65)$$

Hence its determinant can be factored into two functions defined as:

$$\begin{aligned} G^+ &= 1.087 \frac{(s+0.9558)(s+2.0649)(s+0.8518+j0.2437)(s+0.8518-j0.2437)}{(s+0.923)^2(s+1)^2} \\ G^- &= \frac{(s-0.9558)(s-2.0649)(s-0.8518+j0.2437)(s-0.8518-j0.2437)}{(s-0.923)^2(s-1)^2} \end{aligned}$$

The matrices  $\mathbf{N}_1$  and  $\mathbf{A}$  are:

$$\begin{aligned} \mathbf{N}_1 &= \begin{bmatrix} 2 \\ -s^2+1 \\ 2 \\ -s^2+1 \end{bmatrix} \quad (66) \\ \mathbf{A} &= \begin{bmatrix} 0.962 \frac{-s^2+2.367}{-s^2+0.853} & -2 \\ -2 & 1.129 \frac{-s^2+2.144}{-s^2+0.853} \end{bmatrix} \frac{1}{-s^2+1} \quad (67) \end{aligned}$$

By using equation 47:

$$G^+ \mathbf{Y}_1 = \left\{ \frac{1}{G^-} \mathbf{A} \mathbf{N}_1 \right\}^+ + \left\{ \frac{1}{G^-} \mathbf{A} \mathbf{F}_1 \right\}^+ = \begin{bmatrix} \frac{-6.802s-7.503}{(s+1)^2} + \frac{7.373}{s+0.923} + \frac{A_{11}}{s+0.923} + \frac{A_{12}}{s+1} \\ \frac{-6.687s-7.388}{(s+1)^2} + \frac{7.373}{s+0.923} + \frac{A_{21}}{s+1} + \frac{A_{22}}{s+0.923} \end{bmatrix} \quad (68)$$

It is obvious that the right side matrix can be written as:

$$\begin{bmatrix} a_2 s^2 + a_1 s + a_0 \\ (s+1)^2(s+0.923) \\ b_2 s^2 + b_1 s + b_0 \\ (s+1)^2(s+0.923) \end{bmatrix}$$

where  $a_2$ ,  $a_1$ , and  $a_0$  are the new unknown coefficients which are linear combinations of  $A_{11}$  and  $A_{12}$ ; and  $b_2$ ,  $b_1$ , and  $b_0$  are another set of new unknown coefficients which are linear combinations of  $A_{21}$  and  $A_{22}$ . Thus the transfer function vector is:

$$\mathbf{Y}_1 = \begin{bmatrix} \frac{(s+0.923)(a_2 s^2 + a_1 s + a_0)}{(s+0.9558)(s+2.0649)(s+0.8518+j0.2437)(s+0.8518-j0.2437)} \\ \frac{(s+0.923)(b_2 s^2 + b_1 s + b_0)}{(s+0.9558)(s+2.0649)(s+0.8518+j0.2437)(s+0.8518-j0.2437)} \end{bmatrix} \quad (69)$$

Each element is of the form expressed in equation 51. It is observed that quick examination on the right side of equation 68 will yield equation 69 immediately.

By substituting equation 69 into equation 62, six independent equations for these six coefficients will be obtained from evaluating the residues of the left half plane poles. These poles are the four poles of the transfer functions and the one at  $s = -1$ . It should be noted that the two equations obtained by using the same transfer-function's pole in the set of optimum equations 62 will be dependent. However, the two equations associated with the pole at  $s = -1$  are independent. Thus the optimum transfer functions are (equations 70 and 71):

If the input of the first channel alone is used, the optimum filter and its minimum mean-square error are:

$$Y_a(s) = \frac{1 - 0.368e^{-s}}{1 - 0.155e^{-s}} \frac{0.577}{s+1}; \quad (\epsilon_a^2)_{\min} = 0.7364$$

On the other hand, if the input of the second channel only is used, they are:

$$Y_b(s) = \frac{1 - 0.368e^{-s}}{1 - 0.131e^{-s}} \frac{0.6447}{s+1}; \quad (\epsilon_b^2)_{\min} = 0.7072$$

It is clear then that the multipole filter has better performance than the corre-

sponding 2-pole filters at the expense of a more complicated system.

If equation 68 is chosen for the expression of the transfer functions, then the four independent equations for the unknown coefficients can be obtained by evaluating the residues associated with the four poles of the transfer functions either from the first or second equation of the simultaneous equations 62. It should be true, theoretically, that the residues associated with the pole at  $s = -1$  on the left sides of these two systems of equations are 1 and, therefore, are independent of the unknown coefficients. However, because of the computational errors, unfortunately these residues will not be exactly equal to the expected values. On the other hand, the second expressions usually tend to minimize the computational errors which will inevitably occur.

## Conclusions

This paper has shown that with stationary inputs and the assumption of

statistically independent sampling processes, an optimal multipole filter in the Wiener sense can be synthesized.

1. The cross-power spectral densities of the sampled inputs and cross-power spectral densities of the sampled and the cor-



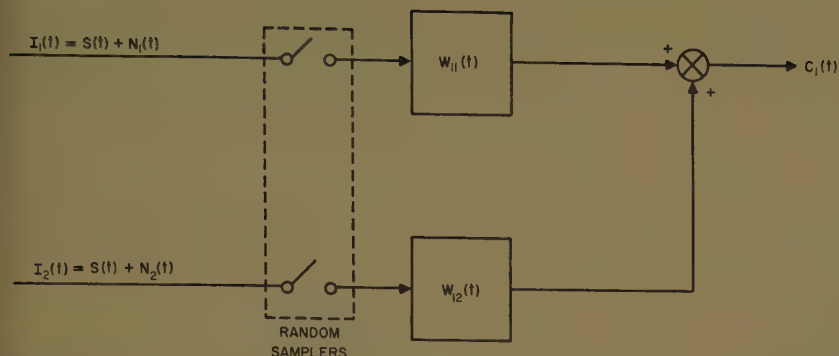


Fig. 8. Example 2, 2+1 pole filter

continuous inputs are unaffected by the sampling processes. The auto-power spectral densities of the sampled inputs are, on the other hand, influenced by the sampling processes.

2. If the sampled spectral densities are not rational functions in  $s$ , they have to be approximated, at least for small  $s$ , by this kind of function. Since the signals usually have power concentrated in the low frequencies, the error introduced by this approximation is negligible.

3. The set of generalized Wiener-Hopf equations is solved by transform method. Factorization of the spectral-density determinant and the method of unde-

termined coefficients are employed in obtaining the optimal transfer functions, which are rational functions in  $s$  and therefore can easily be synthesized with a finite number of lumped linear passive components.

## Appendix. Padé Approximants for $e^{-s}$

$$\frac{1-s}{1} \quad \frac{1-s+\frac{s^2}{2!}}{1}$$

$$\frac{1-\frac{1}{2}s}{1+\frac{1}{2}s} \quad \frac{1-\frac{2}{3}s+\frac{1}{3}\frac{s^2}{2!}}{1+\frac{1}{3}s}$$

$$\frac{1-\frac{1}{3}s}{1+\frac{2}{3}s+\frac{1}{3}\frac{s^2}{2!}} \quad \frac{1-\frac{1}{2}s+\frac{1}{6}\frac{s^2}{2!}}{1+\frac{1}{2}s+\frac{1}{6}\frac{s^2}{2!}}$$

## References

1. ANALYSIS AND SYNTHESIS OF LINEAR SYSTEMS OPERATING ON RANDOMLY SAMPLED DATA, R. E. Kalman. *Doctoral Dissertation*, Department of Electrical Engineering, Columbia University, New York, N. Y., 1957.
2. THE SYNTHESIS OF OPTIMUM RANDOM SAMPLING SYSTEMS, A. R. Bergen. *Technical Report T-T/133*, Electronics Research Laboratories, Columbia University, 1957.
3. EXTRAPOLATION, INTERPOLATION AND SMOOTHING OF STATIONARY TIME SERIES (book), Norbert Wiener. John Wiley & Sons, Inc., New York, N. Y., 1949.
4. ON THE OPTIMUM SYNTHESIS OF MULTIPOLE CONTROL SYSTEMS IN THE WIENER SENSE, H. C. Hsieh, C. T. Leondes. *Transactions, Professional Group on Automatic Control*, Institute of Radio Engineers, New York, N. Y., vol. AC-4, no. 2, Nov. 1959, pp. 16-29.
5. ON THE OPTIMUM SYNTHESIS OF SAMPLED DATA MULTIPOLE FILTERS WITH RANDOM AND NONRANDOM INPUTS, H. C. Hsieh, C. T. Leondes. *Ibid.*, vol. AC-5, no. 3, Aug. 1960, pp. 193-208.

# Sampled-Data Control Systems with Transport Lag by Mitrović's Algebraic Method

DRAGOSLAV ŠILJAK  
NONMEMBER AIEE

TRANSPORT LAG in sampled-data control systems influences their stability. If this influence is studied as the function of transport lag and other system parameters, the design procedure becomes cumbersome. This paper, which applies Mitrović's algebraic method,<sup>1-3</sup> presents a new way to analyze and synthesize various sampled-data control systems such as the cyclic variable-rate, multirate, and those with nonsynchronized samplers. Other related problems are greatly facilitated by this method also.

Paper 61-743, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE Summer General Meeting, Ithaca, N. Y., June 18-23, 1961. Manuscript submitted February 21, 1961; made available for printing April 25, 1961.

DRAGOSLAV ŠILJAK is with the University of Belgrade, Belgrade, Yugoslavia.

Making use of flat-top-pulse or trapezoid-pulse approximations,<sup>4</sup> the proposed design procedure can be extended to cover sampled-data control systems with finite sampling duration. Moreover, this technique can deal effectively with problems appearing in sampled-data systems in which hidden instability<sup>5,6</sup> occurs.

A simple single-loop system with transport lag  $T_t$  and the conventional linear part characterized by a transfer function  $G(s)$ , shown in Fig. 1, will be considered representative in the analysis. Procedure and conclusions outlined herein are applicable, with only slight modifications, to other types mentioned.

Before analyzing the system, a necessary preliminary is to determine the  $z$ -transfer function  $G(z, T_t)$ , corresponding to  $e^{-T_t s} G(s)$ . Since the method

to be used is based upon study of the system-characteristic equation having the algebraic form

$$a_n z^n + a_{n-1} z^{n-1} + \dots + a_1 z + a_0 = 0 \quad (1)$$

$G(z, T_t)$  should be a rational function in  $z$ .  $T_t$  is not generally a multiple integer of the sampling period  $T$ . Hence, to obtain  $G(z, T_t)$  in the required form,  $T_t$  must be expressed as

$$T_t = kT - mT \quad (2)$$

where  $k$  is the upper integer of  $T_t/T$  and  $m$  is a fraction ( $0 < m \leq 1$ ). Finding the modified  $z$ -transform  $G(z, m)$  corresponding to  $G(s)$  and substituting  $m$  with  $k - T_t/T$  from equation 2, the required form of function  $G(z, T_t)$  is

$$G(z, T_t) = z^{-(k-1)} G(z, m) |_{m=k-T_t/T} = z^{-(k-1)} G(z, k - T_t/T) \quad (3)$$

Using equation 3, system behavior is determined by an algebraic equation of form 1, and Mitrović's algebraic method is applied. For convenience, this method is stated briefly.

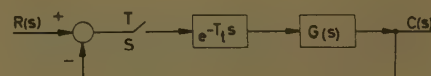


Fig. 1. Block diagram of sampled-data control system with transport lag

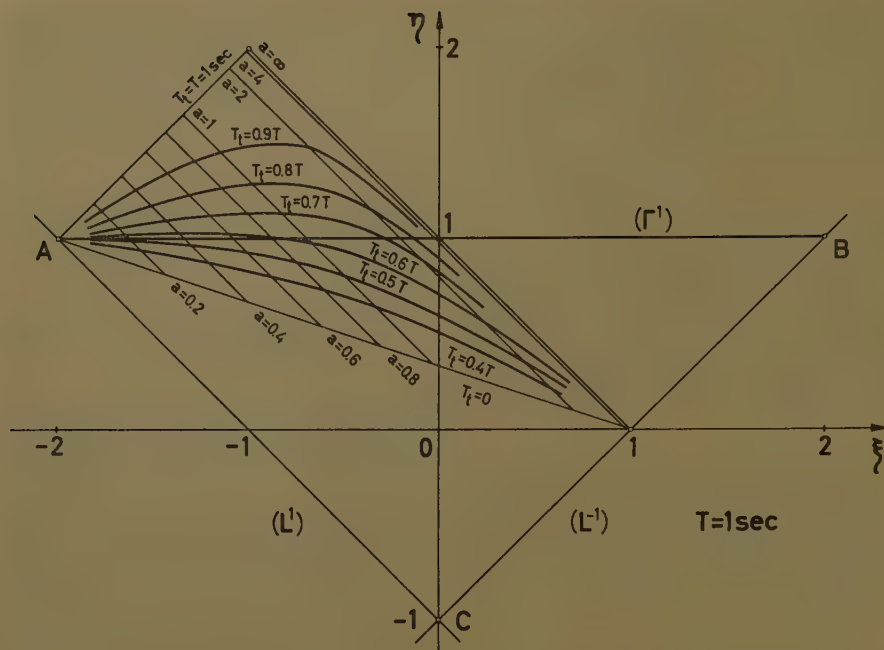


Fig. 2. Correlation between system parameters  $a$  and  $T_t$

Considering coefficients  $a_1$  and  $a_0$  in equation 1 as variables  $\xi$  and  $\eta$ , and using relationships

$$z = -\omega_n \zeta + j\omega_n \sqrt{1-\zeta^2} = e^{sT} = e^{(-\omega_n \zeta + j\omega_n \sqrt{1-\zeta^2})T} \quad (4)$$

where  $\omega_n$  is undamped natural frequency and  $\zeta$  is the relative damping coefficient, will get

$$\begin{aligned} \xi &= a_2 \phi_2(-\cos \omega_n T \sqrt{1-\zeta^2}) e^{-\omega_n \zeta T} + \\ &\quad a_3 \phi_3(-\cos \omega_n T \sqrt{1-\zeta^2}) e^{-2\omega_n \zeta T} + \dots \\ \eta &= -e^{-2\omega_n \zeta T} [a_2 \phi_1(-\cos \omega_n T \sqrt{1-\zeta^2}) + \\ &\quad a_3 \phi_2(-\cos \omega_n T \sqrt{1-\zeta^2}) e^{-\omega_n \zeta T} + \dots] \end{aligned} \quad (5)$$

These equations represent the loci of points in the  $0\xi\eta$ -plane corresponding to roots with settling time, damping coefficient, or undamped natural frequency being constant, depending on which

variable among them is considered as constant. Plotting of these characteristic curves offers no particular difficulty since the table of functions  $e^{-x}$  and Table I in reference 3 of the functions  $\phi_i(-\cos x)$  may be used.

The essential point after plotting these curves in the  $0\xi\eta$ -plane, using equations 1, 4, and 5, is to study the position of working point  $M(a_1; a_0)$  in the same plane. For sampled-data control systems with transport lag, one great advantage of the described method is its ability to provide correlation between over-all system performance and time  $T_t$ .

### Absolute Stability

Attainment of absolute stability is now sought. The first step is to extend

the analysis procedure to learn how transport lag and other parameters will affect the stability of the system. A simple example will be presented but the same principles can be used with minor modifications, for other configurations.

The representative system shown in Fig. 1 has the transfer function

$$G(s) = \frac{Ka}{s(s+a)}$$

Applying equation 3, the transform  $G(z, T_t)$  corresponding to this system is

$$G(z, T_t) = z^{-(k-1)} K \times \left[ \frac{1}{z-1} - \frac{1}{z-e^{-aT}} e^{-a(kT-T_t)} \right] \quad (7)$$

When the time  $T_t$  is less than one sampling period  $T$  or

$$k=1 \quad (8)$$

By substituting this value in equation 7 the roots of characteristic equation  $1+G(z, T_t)=0$  are determined by

$$z^2 + (K[1-e^{-a(T-T_t)}] - [1+e^{-aT}])z + Ke^{-a(T-T_t)} + e^{-aT}(1-K) = 0 \quad (9)$$

As known, the system is stable if all roots of the characteristic equation in  $z$  lie within the unit circle around the  $z$ -plane's origin.

The foregoing statement finds expression in terms of analysis in the  $0\xi\eta$ -plane if the limits of the stable region in that plane are determined by mapping the unit circle ( $\omega_s=1$ ) from the  $z$ -plane into the  $0\xi\eta$ -plane.

The stable region in the  $0\xi\eta$ -plane of the representative system is bounded by

1. The locus of points  $(\Gamma^1)$  corresponding to the roots with constant  $\omega_0=1$ , which is the limit curve for complex roots.

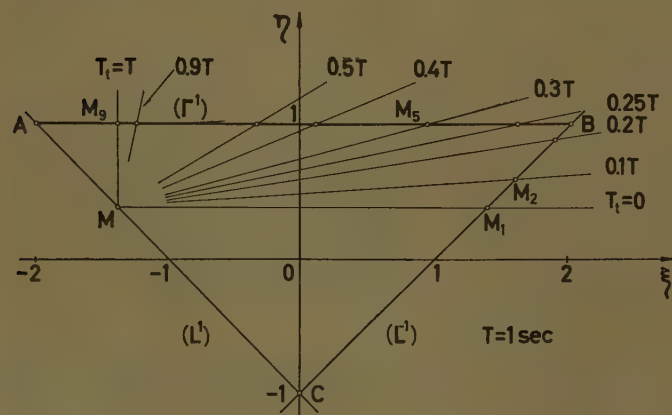


Fig. 3. Determination of gain  $K_{\max}$  for  $0 \leq T_t \leq T = 1$  sec

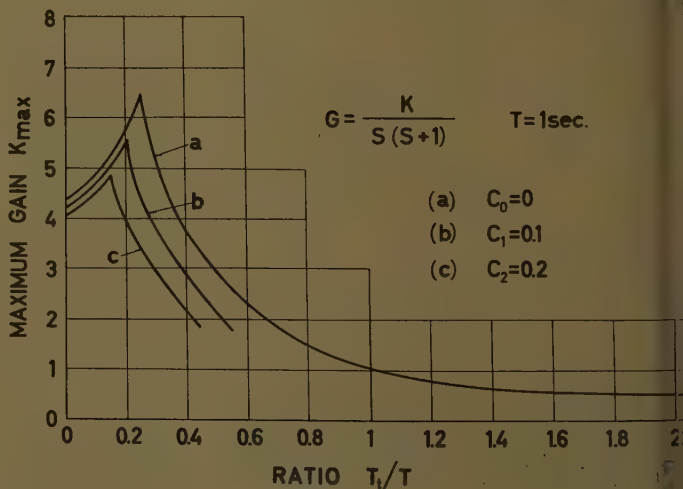


Fig. 4. Diagram relating gain  $K_{\max}$  to ratio  $T_t/T$  for various values of constant  $c$



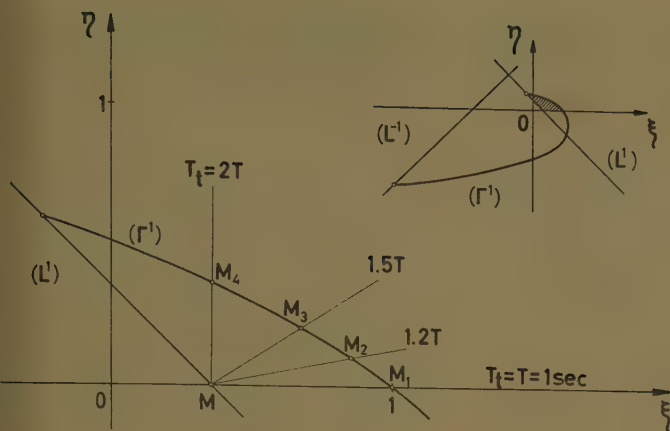


Fig. 5. Determination of gain  $K_{\max}$  for  $T \leq T_i \leq 2T$ ,  $T = 1$  sec

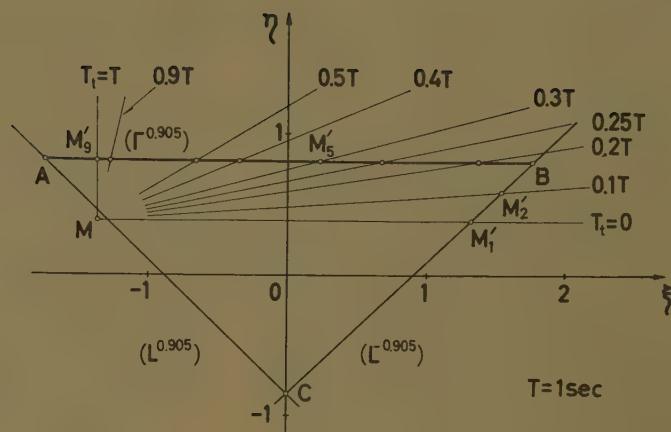


Fig. 6. Determination of gain  $K_{\max}$  for  $c_1 = 0.1$

Since  
 $-\omega_n \xi^2 T = \omega_z = 1$

and

$$\xi_z = -\cos \omega_n T \sqrt{1 - \xi^2} \quad (10)$$

the equations 5 become

$$\begin{aligned} &= \phi_2(\xi_z) \\ &= -\phi_1(\xi_z) \quad (\Gamma^1) \end{aligned} \quad (11)$$

Because

$$\phi_1(\xi_z) = -1$$

$$\phi_2(\xi_z) = 2\xi_z$$

see reference 3), equations 11 of the locus of points, corresponding to the roots with constant  $\omega_z = 1$ , obtain the form

$$\begin{aligned} &= 2\xi_z \\ &= 1 \quad (\Gamma^1) \end{aligned} \quad (12)$$

where  $\xi$  is the variable parameter. As  $\xi$  varies only between  $-1$  and  $+1$ , the equations 12 represent the straight line between  $A$  and  $B$  in Fig. 2.

The locus of points  $(L^1)$ , corresponding to positive unit roots ( $z = +1$ ), which is the limit curve for positive roots.

The locus  $(L^1)$  may be written simply by letting  $z = +1$  in equation 1 and by taking coefficients  $a_1$  and  $a_0$  as variables  $\xi$  and  $\eta$ . So, in the case of equation 9

$$\eta = -\xi - 1 \quad (L^1) \quad (13)$$

which represents line  $(L^1)$  in Fig. 2.

The locus of points  $(L^{-1})$ , corresponding to negative unit roots ( $z = -1$ ), which is the limit curve for negative roots.

The locus  $(L^{-1})$  is written in a form similar to that of  $(L^1)$

$$\eta = \xi - 1 \quad (L^{-1}) \quad (14)$$

which represents line  $(L^{-1})$  in Fig. 2.

Now, the condition of absolute stability, in terms of  $0\xi\eta$ -plane analysis, is

$$\begin{aligned} &K[1 - e^{-a(T-T_i)}] - [1 + e^{-aT}]; \\ &Ke^{-a(T-T_i)} + e^{-aT}(1-K)) \end{aligned}$$

and is inside the stable region bounded

by the curves  $(\Gamma^1)$ ,  $(L^1)$ , and  $(L^{-1})$  in Fig. 2.

Note that co-ordinates of working point  $M$  are functions of system parameters  $a$ ,  $K$ ,  $T$ , and  $T_i$ . Thus, by plotting in the same  $0\xi\eta$ -plane, the stable region and loci of working points  $M$  for different values of system parameters, the influence of each parameter variation upon system stability is easily determined.

When gain  $K$  and sampling period  $T$  are held constant, correlation between parameter  $a$  and transport lag  $T_i$  is found by first assuming that  $K = 2$  and  $T = 1$  sec (second).

Substituting these values of  $K$  and  $T$  in co-ordinates of working point  $M$  gives

$$\begin{aligned} \xi &= 1 - 2e^{-a(1-T_i)} - e^{-a} \\ \eta &= 2e^{-a(1-T_i)} - e^{-a} \end{aligned} \quad (15)$$

For various values of  $a$  and  $T_i$ , equations 15 represent the net of two families of curves as shown in Fig. 2. From the part of this net within the triangular  $ABC$  or stable region, interpolation can be employed in finding any pair of  $a$  and  $T_i$  values which makes the system stable.

Attention also may be focused on correlating gain  $K$  and transport lag  $T_i$ . For example, suppose  $a = 1$  and  $T = 1$  sec, the relation between maximum gain for stability limit  $K_{\max}$  and transport lag  $T_i$  is to be found. By inserting the given values for  $a$  and  $T$  into equations 15, the co-ordinates of working point  $M$  become

$$\begin{aligned} \xi &= [1 - e^{-(1-T_i)}]K - 1.368 \\ \eta &= [e^{-(1-T_i)} - 0.368]K + 0.368 \end{aligned} \quad (16)$$

Evidently, if  $T_i$  is held constant and  $K$  is allowed to vary, the stable region remains unchanged, but working point  $M$  moves along the straight line in the  $0\xi\eta$ -plane determined by equations 16. For various values of  $T_i$ , equations 16

represent the family of straight lines plotted in Fig. 3. Now, from the intersection points  $(M_1, M_2, \dots, M_9)$  of these straight lines and the lines limiting the stable region, by using one of equations 16, the corresponding values of gain  $K_{\max}$  are determined graphically.

In the simple example, analytical solution is also possible. It consists of determining the intersection points  $(M_1, M_2, \dots, M_9)$  with the help of equations 12, 14, and 16. When the co-ordinates  $\xi$  and  $\eta$  of these points are obtained, the corresponding equation in equations 16 will yield the gain values  $K_{\max}$ . In doing so, the first part of the curve  $a$  in Fig. 4, which corresponds to the case when  $T_i$  is less than  $T$ , can be obtained.

So far, the analysis has been carried out for the case  $0 \leq T_i \leq T$ . The second part of curve  $a$ , which corresponds to the second sampling period, is given by first inserting into equation 7

$$k = 2 \quad (17)$$

Now, for  $T = 1$  sec and  $a = 1$ , the corresponding characteristic equation is

$$z^3 - 1.368z^2 + (K[1 - e^{-(2-T_i)}] + 0.368)z + K[e^{-(2-T_i)} - 0.368] = 0 \quad (18)$$

Using the same procedure as for system equation 9, the curves which define the stable region, corresponding to system equation 18, are determined by

$$\begin{aligned} \xi &= -1.368\phi_2(\xi_z) + \phi_3(\xi_z) \\ \eta &= 1.368\phi_1(\xi_z) - \phi_2(\xi_z) \quad (\Gamma^1) \\ \eta &= -\xi + 0.368 \quad (L^1) \\ \eta &= \xi + 2.368 \quad (L^{-1}) \end{aligned} \quad (19)$$

The loci of points  $(\Gamma^1)$ ,  $(L^1)$ , and  $(L^{-1})$  have the form shown in the upper right corner in Fig. 5.

From system equation 18, the co-ordinates of working point  $M$  are

$$\begin{aligned} \xi &= [1 - e^{-(2-T_i)}]K + 0.368 \\ \eta &= [e^{-(2-T_i)} - 0.368]K \end{aligned} \quad (20)$$

Both the enlarged part of the stable region and the straight lines defined by equations 20 for different values of  $T_i$  are shown in Fig. 5. Applying the same procedure used for the first sampling period, the second part of curve  $a$  is plotted in Fig. 4. The computation necessary to produce such diagrams for more complicated configurations can be performed by the proposed design procedure with relative ease and in a systematic manner. In designing sampled-data control systems, where transport lag is introduced to provide stabilization, the proposed procedure is extremely useful.

Where stabilization and compensation by conventional means are difficult, sampling can be used often—if the sampling period  $T$  is properly chosen—to improve stability. Thus, when applying the proposed technique, an optimum value of the sampling rate can be designed without creating new problems. Only sampling period  $T$  and gain  $K$  need to be considered as variable parameters. By investigating positions of working point  $M$ , as in the foregoing examples, suitable magnitude of the sampling rate can be determined easily.

## Settling Time

In sampled-data control systems, the prescribed maximum settling time often appears as a requirement. This usually signifies that the adjustable parameters of a given system must be set so that the system will settle to an essentially steady-state value in the required time after application of a step-input function.

Settling time will not exceed the prescribed maximum value if all poles of the system-transfer function have a real part greater than some corresponding constant value  $c$ ; that is, if

$$\omega_n t^* \geq c$$

Because of equation 4, the foregoing condition means that all zeros of the characteristic equation in  $z$  should lie within the circle whose radius is

$$\omega_z = e^{-\omega_n t^*} = e^{-cT} \quad (21)$$

Since this condition is so similar to that of absolute stability ( $\omega_z = 1$ ), the design procedure which determines a system with a prescribed settling time is essentially that used in the preceding section. The only difference is that the radius of the circle in the  $z$ -plane is not unity.

To illustrate the design procedure, consider again the system shown in Fig. 1 with the transfer function  $G(s)$  given in equation 6. Suppose parameter

values are again  $a=1$  and  $T=1$  sec, and values are required for maximum gain  $K_{\max}$  and transport lag  $T_i$ , which will result in a settling time corresponding to  $c_1=0.1$ .

Since the system will be investigated for  $T_i$  values that are less than the sampling period  $T$ —that is, for  $k=1$ —the corresponding characteristic equation is 9, into which is inserted  $a=1$  and  $T=1$  sec, giving the following equation

$$z^2 + (K[1 - e^{-(1-T)}] - 1.368)z + K(e^{-(1-T)} - 0.368) + 0.368 = 0 \quad (22)$$

Settling time will be less than the prescribed maximum value determined by  $c_1=0.1$  if the roots of equation 22 lie inside the circle whose radius is

$$\omega_z = e^{-c_1 T} = 0.905 \quad (23)$$

Limiting curves in the  $0\xi\eta$ -plane, corresponding to this circle, are obtained by the same method described in the foregoing section. Curves  $(\Gamma^{0.905})$ ,  $(L^{0.905})$ , and  $(L^{-0.905})$  consequently are defined by

$$\xi = a_2 \phi_2(\xi_z) e^{-c_1 T} = 1.81 \xi_z$$

$$\eta = -a_2 \phi_1(\xi_z) e^{-2c_1 T} = 0.82 \quad (\Gamma^{0.905})$$

$$\eta = e^{-c_1 T} \xi - a_2 e^{-2c_1 T} = 0.905 \xi - 0.82 \quad (L^{-0.905})$$

$$\eta = -e^{-c_1 T} \xi - a_2 e^{-2c_1 T} = -0.905 \xi - 0.82 \quad (L^{0.905}) \quad (24)$$

These curves are shown in Fig. 6.

Now, in terms of  $0\xi\eta$ -plane analysis, the condition for maximum settling time may be formulated. Settling time will be less than the prescribed value corresponding to  $c_1=0.1$ , if the working point

$$M(K[1 - e^{-(1-T)}] - 1.368; K[e^{-(1-T)} - 0.368] + 0.368)$$

is bounded by curves determined in equations 24.

Since the co-ordinates of point  $M$  are the same as in equations 16, the system of straight lines plotted in Fig. 6 is the same as was plotted in Fig. 3. The same graphical procedure used in the foregoing section will produce curve  $b$  in Fig. 4 from points  $M'_1, M'_2, \dots, M'_9$ . In like manner, if the required maximum settling time corresponds to  $c_2=0.2$ , curve  $c$  in Fig. 4 can be plotted.

Bearing in mind that curve  $a$  corresponds to  $c_0=0$ —that is,  $\omega_z=1$ —the tendency for gain  $K_{\max}$  to change realized, the amount depending on the magnitude of the maximum settling time prescribed and the transport lag. By using interpolation, diagrams Fig. 4 are useful in determining the pairs of parameter values  $K$  and  $T$  which will yield the desired settling time. Since similar diagrams are obtainable if the relative damping is specified (considering in equations 5 relative damping coefficient  $\zeta$  as constant),<sup>3</sup> the proposed method is applicable in designing systems with desired transient performance. For systems which hidden oscillations<sup>5,6</sup> can take place, the results obtained in the design of transient behavior by this procedure should be taken with some caution.

## Nonsynchronized Samplers

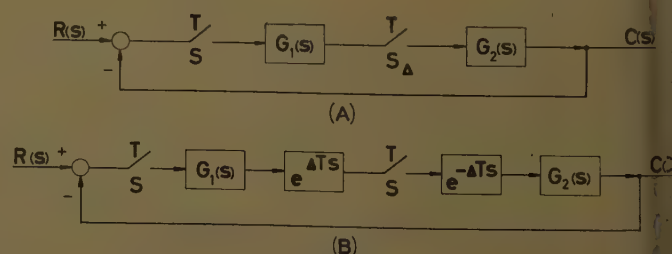
Another application of the proposed method is to systems containing two or more samplers operating with the same sampling rates, but not synchronized in phase, commonly known as sampled-data systems with nonsynchronized samplers.<sup>4</sup>

The block diagram in Fig. 7(A) represents the simplest type, having only two samplers  $S$  and  $S_\Delta$  with the same sampling period  $T$ . However,  $S_\Delta$  "slips" behind  $S$  by  $\Delta T$  sec, where  $\Delta$  is a fraction, and may be referred to as the slip factor. Based on the concept of equivalent samplers,<sup>4</sup>  $S_\Delta$  may be represented by a basic sampler  $S$ , preceded by an advance of  $\Delta T$  sec and followed by  $\Delta T$  sec delay as in Fig. 7(B).

With reference to Fig. 7(B), the open-loop  $z$ -transfer function is  $G(z, \Delta) = \mathfrak{z}[e^{\Delta T s} G_1(s)] \mathfrak{z}[e^{-\Delta T s} G_2(s)]$  where the notation  $\mathfrak{z}$  indicates the  $z$ -transform corresponding to the function which it prefixes. When  $G_1(s)$  and  $G_2(s)$  are rational functions in  $s$ , it follows from equation 3 and the definition of the modified  $z$ -transform that  $\mathfrak{z}[e^{\Delta T s} G_1(s)] = z G_1(z, m)|_{m=\Delta} = z G_1(z, \Delta)$   $\mathfrak{z}[e^{-\Delta T s} G_2(s)] = G_2(z, m)|_{m=1-\Delta} = G_2(z, 1-\Delta)$

Where the roots of the corresponding characteristic equation are located within the unit circle, the system is stable.

Fig. 7. A—Block diagram of basic nonsynchronized sampled-data system. B—Equivalent block diagram of system A





etermine stability of the system represented in Fig. 7(A); see equation 27, which shows the influence of the slip factor.

$$+zG_1(z, \Delta)G_2(z, 1-\Delta)=0 \quad (27)$$

When the values of  $\Delta$  fall within a certain range, the nonsynchronized sampling usually provides a stabilizing effect. In a distinct class of sampled-data systems, nonsynchronized sampling is introduced deliberately to improve stability. In designing such systems, the slip factor's optimum value frequently must be determined to find the highest stability limit compatible with maximum allowable gain. To illustrate use of the proposed technique, a simple numerical example is presented, based on the block diagram, Fig. 7(A). Transfer functions  $G_1(s)$  and  $G_2(s)$  are

$$G_1(s)=\frac{K(1-e^{-sT})}{s(s+1)} \quad G_2(s)=\frac{1}{s+0.5} \quad (28)$$

For both samplers, the sampling period is  $T=1$  sec. The effect of slip factor  $\Delta$  on system stability is now studied.

Making use of equations 25, 26, and 27, the corresponding characteristic equation is

$$z^2 - (0.97 - K[e^{-0.5(1-\Delta)} - e^{-0.5(1+\Delta)}])z + K[e^{-0.5(1+\Delta)} - 0.368e^{-0.5(1-\Delta)}] + 0.22 = 0 \quad (29)$$

Since the shape of characteristic curves  $\Gamma^1$ ,  $(L^1)$ , and  $(L^{-1})$  does not depend on values of coefficients  $a_1$  and  $a_0$ , equation 1, the stable region corresponding to equation 29 is the same as that shown in Figs. 2 and 3, which corresponds to the characteristic equation 9.

Co-ordinates of working point  $M$  are

$$\begin{aligned} &= [e^{-0.5(1-\Delta)} - e^{-0.5(1+\Delta)}]K - 0.97 \\ &= [e^{-0.5(1+\Delta)} - 0.368e^{-0.5(1-\Delta)}]K + 0.22 \end{aligned} \quad (30)$$

For various values of slip factor  $\Delta$  and gain constant  $K$ , a system of straight lines similar to that shown in Fig. 3 can be constructed. Then, using the same

graphical procedure as was applied to Fig. 3, diagram (a) in Fig. 8 is plotted. Now, obviously, any value of slip factor  $\Delta$ , which lies between 0 and 1, will improve system stability. The optimum value, procured from this diagram, is 0.76.

If transfer function  $G_1(s)$  has the form

$$G_1(s)=\frac{Ke^{-sT}}{(s+1)^2+1} \quad (31)$$

and transfer function  $G_2(s)$  remains unchanged, the corresponding characteristic equation is

$$\begin{aligned} z^3 - 0.753z^2 + 0.607Ke^{-1.5\Delta} \sin \Delta + \\ 0.106z + 0.082Ke^{-1.5\Delta} \sin(1-\Delta) - \\ 0.01 = 0 \end{aligned} \quad (32)$$

In like manner, by applying the proposed technique, diagram (b) in Fig. 8 is evolved. In contrast to the preceding case, where nonsynchronized sampling provided a stabilizing effect, here the introduction of any slip factor value decreases the stability margin.

Plotting of diagrams that correspond to those in Fig. 8 reveals whether deterioration or improvement will result from nonsynchronized sampling in a given system, and shows the optimum value of slip factor  $\Delta$ . The computation is relatively easy. Diagrams relating  $\Delta$  to other parameters can also be plotted. They prove quite useful when designing sampled-data control systems where nonsynchronized sampling occurs inherently and where it is required for stabilizing effects.

## Conclusions

The proposed design procedure is a graphical means of analyzing and synthesizing sampled-data control systems in which transport lag occurs inherently or is introduced deliberately to provide stabilization. It consists of plotting some characteristic curves in the  $0\eta\zeta$ -plane and then investigating the positions of working point  $M$  as the function of system parameters. In most cases, all

criteria<sup>5</sup> (which are a further simplification of Routh-Hurwitz criteria) for continuous systems. These will be stated now for the second- and third-order case and then used to find boundaries of the stable region in the  $0\eta\zeta$ -plane as well as values for  $K_{\max}$ . Criteria for the general  $n$ th-order case are contained in reference 4.

### STABILITY TEST, SECOND ORDER CASE<sup>4</sup>

where  $n=2$ ,  $F(z)=a_2z^2+a_1z+a_0$ ,  $a_2>0$

$$\begin{aligned} |a_0| < a_2 \\ |a_0+a_2| > |a_1| \text{ or } a_0+a_1+a_2 > 0 \\ a_0-a_1+a_2 > 0 \end{aligned} \quad (33)$$

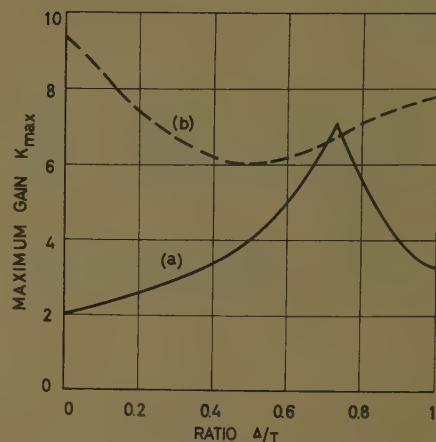


Fig. 8. Diagrams relating gain  $K_{\max}$  to ratio  $\Delta/T$  for two different systems

corresponding curves need not be plotted, and often the determination of only part of a curve is necessary (as shown in Fig. 5). Computations can be performed with ease.

The technique is applicable to many varied problems that appear in sampled-data systems with two or more samplers. Analysis and synthesis of such systems, having finite sampling duration, can also be performed in a straightforward manner.

## References

1. GRAPHICAL ANALYSIS AND SYNTHESIS OF FEEDBACK CONTROL SYSTEMS, I—THEORY AND ANALYSIS, Dušan Mitrović. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 77, 1958 (Jan. 1959 section), pp. 476-87.
2. GRAPHICAL ANALYSIS AND SYNTHESIS OF FEEDBACK CONTROL SYSTEMS, II—SYNTHESIS, Dušan Mitrović. *Ibid.*, pp. 487-96.
3. GRAPHICAL ANALYSIS AND SYNTHESIS OF FEEDBACK CONTROL SYSTEMS, III—SAMPLED-DATA FEEDBACK CONTROL SYSTEMS, Dušan Mitrović. *Ibid.*, pp. 497-503.
4. DIGITAL AND SAMPLED-DATA CONTROL SYSTEMS (book), J. T. Tou. McGraw-Hill Book Company, Inc., New York, N. Y., 1959.
5. SAMPLED-DATA CONTROL SYSTEMS (book), Eliahu I. Jury. John Wiley & Sons, Inc., New York, N. Y., 1958.
6. HIDDEN OSCILLATIONS IN SAMPLED-DATA CONTROL SYSTEMS, Eliahu I. Jury. *Ibid.*, vol. 75, 1956 (Jan. 1957 section), pp. 391-95.
7. THE ANALYSIS OF SAMPLED-DATA SYSTEMS, J. R. Ragazzini, L. A. Zadeh. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 71, Nov. 1952, pp. 225-34.

### STABILITY TEST, THIRD ORDER CASE<sup>4</sup>

where  $n=3$ ,  $F(z)=a_3z^3+a_2z^2+a_1z+a_0$ ,  $a_3>0$

$$\begin{aligned} |a_0| < a_3 \\ a_0^2 - a_3^2 < a_0a_2 - a_1a_3 \\ a_0 + a_1 + a_2 + a_3 > 0 \\ a_0 - a_1 + a_2 - a_3 < 0 \end{aligned} \quad (34)$$

### EXAMPLE 1

Considering characteristic equation 9 of the paper, recognizing that  $a_0=\eta$ ,  $a_1=\xi$  and using equations 33 readily produces

$$\begin{aligned} |\eta| &< 1 & (35) \\ \eta + 1 + \xi &> 0 & (36) \\ \eta - \xi + 1 &> 0 & (37) \end{aligned}$$

The three inequalities give the stability region of Fig. 3 of this paper.

#### EXAMPLE 2

Considering equation 18 of the paper, which is a polynomial in  $s$  of third degree, and using equation 34 will give the following:

$$\begin{aligned} |\eta| &< 1 & (38) \\ \eta^3 - 1 &< -1.368\eta - \xi & (39) \\ \eta + \xi - 1.368 + 1 &> 0 & (40) \\ \eta - \xi - 1.368 - 1 &< 0 & (41) \end{aligned}$$

Careful interpretation of these inequalities will reveal the stability region shown in Fig. 5 of the paper.

#### EXAMPLE 3

Consider an example of synthesis to find the  $K_{\max}$  required for a given system to remain stable. The stability test of equation 32 of the paper is

$$|0.082Ke^{-1.5\Delta} \sin(1-\Delta) - 0.01| < 1 \quad (42)$$

$$\begin{aligned} (0.082Ke^{-1.5\Delta} \sin(1-\Delta) - 0.01)^2 - & \\ 1 < -0.753(0.082Ke^{-1.5\Delta} \sin(1-\Delta) - & \\ 0.01) - (0.607Ke^{-1.5\Delta} \sin \Delta + 0.106) & (43) \\ 0.082Ke^{-1.5\Delta} \sin(1-\Delta) - 0.01 + & \\ 0.607Ke^{-1.5\Delta} \sin \Delta + 0.106 - 0.753 + 1 > 0 & (44) \\ 0.082Ke^{-1.5\Delta} \sin(1-\Delta) - 0.01 - & \\ 0.607Ke^{-1.5\Delta} \sin \Delta - 0.106 - 0.753 - 1 < 0 & (45) \end{aligned}$$

For a given value of  $\Delta$ , relations 42 through 45 determine the allowable values of  $K$ . Thus,  $K_{\max}$  is readily obtained from the inequalities. For example, take  $\Delta=1$  corresponding to  $\Delta/T=1$  in Fig. 8, where equations 42-45 yield relations 46-49,

$$\begin{aligned} 0.01 &< 1 & (46) \\ \text{This gives no information on } K. & & \\ K &< 7.93 & (47) \\ K &> -3.07 & (48) \\ K &> -14.7 & (49) \end{aligned}$$

Since only positive values of  $K$  are considered,  $K_{\max}=7.93$ , which agrees with the value corresponding to  $\Delta/T=1$  in curve (b) of Fig. 8. Similarly,  $K_{\max}$  could also be obtained for other values of  $\Delta$ . This method could be applied also to examples contained in reference 5 of the

paper. Stability tests in reference 4 will give the criteria for absolute stability. In the case of instability, it also indicates the number of roots which lie inside or outside the unit circle.

While the paper gives a useful graphical design technique, the present discussion outlines an alternate method which achieves the same results by using Jury's simplified stability criteria.<sup>4</sup> Thus, the design engineer has a choice.

#### REFERENCES

1. See references 1, 2, and 3 of the paper.
2. NOTES ON THE STABILITY CRITERION FOR LINEAR DISCRETE SYSTEMS, Eliahu I. Jury, B. H. Bhargava. *Transactions, Professional Group on Automatic Control, Institute of Radio Engineers*, New York, N. Y., vol. AC-6, no. 1, Feb. 1961, pp. 89-90.
3. ADDITIONS TO STABILITY CRITERION FOR LINEAR DISCRETE SYSTEMS, Eliahu I. Jury. *Ibid.*, vol. AC-6, no. 1, Sept. 1961, pp. 342-43.
4. A SIMPLIFIED STABILITY CRITERION FOR LINEAR DISCRETE SYSTEMS, Eliahu I. Jury. *Report no. 373, Electronics Research Laboratory, University of California, Berkeley, Calif.*, Series no. 60, June 1961.
5. *Theory of Matrices* (book), F. R. Gantmacher, Chelsea Publishing Co., New York, N. Y., vol. 1, 1959, pp. 221-27.

Dragoslav Šiljak: I read with interest the discussion by M. Anantha Pai and do not consider additional comment necessary.

# Mathematical Models for Time-Domain Design of Electrohydraulic Servomechanisms

P. K. C. WANG  
ASSOCIATE MEMBER AIEE

IN THE ANALYTICAL DESIGN of high-speed electrohydraulic servomechanisms, an accurate mathematical model of the valve-actuator mechanism is generally helpful in determining the controller and system configuration required to achieve the desired over-all performance. During past years, various forms of electrohydraulic valves and actuator mechanisms have been realized. Detailed descriptions of these devices can be found in a book by Blackburn et al.<sup>1</sup> Most previous works on the analysis of these systems have resorted to lineariza-

tion techniques.<sup>2,3</sup> Recently, Butler and Turnbull presented detailed analysis of simplified nonlinear models.<sup>4,5</sup> Also, Zaborzsky and Harrington derived describing functions for both single and two-stage electrohydraulic valves, which permit analytical design in the frequency-domain.<sup>6-9</sup> The limitations of linearized models are discussed by Rausch.<sup>10</sup>

The objective of this paper is to derive mathematical models which will

represent most conventional electrohydraulic valve-controlled actuators with sufficient accuracy for time-domain design and to examine the system from a general viewpoint so as to reveal some of its salient features. Emphasis is placed upon derivation of time-domain trajectories for various types of inputs. The results form a basis for further investigations on sampled-data and time-optimal control of electrohydraulic systems.

## System Description

The system under consideration is shown in Fig. 1. The functional parts consist of a 4-way valve with a generalized valve-spool driver and a piston-type actuator rigidly attached to the load. The valve-spool driver is usually an electric torque-motor which is either connected directly to the spool or indirectly through additional stages of hydraulic

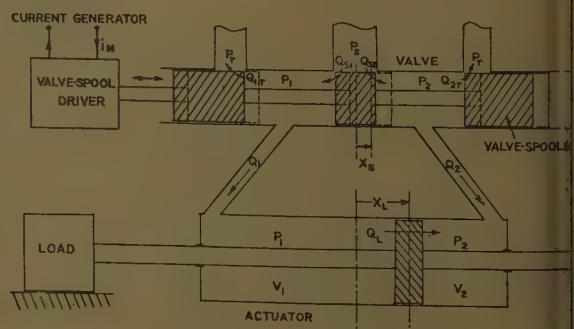


Fig. 1. A typical electrohydraulic valve-controlled actuator

Paper 61-819, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE Summer General Meeting, Ithaca, N. Y., June 18-23, 1961. Manuscript submitted October 14, 1960; made available for printing April 21, 1961.

P. K. C. WANG is with the International Business Machines Corporation, San Jose, Calif.

The author wishes to thank J. T. Ma and J. O. Hildebrand for their many helpful discussions.



plification. The electric torque-motor assumed to be driven by an ideal current generator so as to eliminate the effect equivalent motor armature series inductance. Also, the actuator load consists of mass  $M_L$  and a general friction load

The valve-spool displacement  $x_s(\tau)$  is measured from a neutral position where the spool is symmetrically centered with respect to the valve ports. The actuator inlet displacement,  $x_L(\tau)$ , is also defined from a fixed reference corresponding to the middle point of the actuator cylinder. Both  $x_s(\tau)$  and  $x_L(\tau)$  are chosen to be positive from left to right.

## Basic Assumptions

The valve is symmetrical (that is, the inlet and outlet ports are identical).

The oil in the valve cylinder is incompressible and the flow is steady.

The volumetric flow  $Q$  through the valve ports can be related to the pressure differential  $\Delta P$  and port area  $A_p$  instantaneously by

$$Q = C_p A_p \operatorname{sgn} \Delta P \sqrt{|\Delta P|} \quad (1)$$

where  $C_p$  is an average discharge proportionality constant.

The valve is located near the actuator so that pressure drops and propagation time lag due to connecting lines are negligible. Also, the pressures are uniform inside the actuator cylinder volumes  $V_1$  and  $V_2$ .

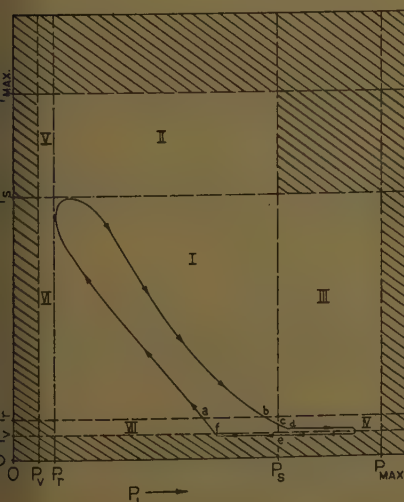
The supply pressure is maintained at constant value  $P_s$ .

Effects of temperature variations and aging on the oil properties are negligible.

Further assumptions will be mentioned, where applicable, in the ensuing discussion.

## Model for Large Motions

The equation of motion for the valve spool is



$$M_s \ddot{x}_s(\tau) + F_s[\dot{x}_s(\tau)] + K_s x_s(\tau) = F_R + F_T[i_M(\tau)] \quad (2)$$

where  $M_s$  is the spool mass,  $F_s$  is a friction force function depending upon the spool velocity,  $K_s$  is a valve spring constant, and  $F_T$  is assumed to be an amplitude-sensitive nonlinear function of the actual torque-motor current  $i_M(\tau)$ . In particular,  $F_T$  may be a hysteresis nonlinearity caused by magnetic hysteresis in the motor.  $F_R$  is the axial component of Bernoulli reaction force which is induced by a change in the momentum of oil discharging through the valve ports. Lee and Blackburn have shown that  $F_R$  can be decomposed into steady-state and transient components.<sup>11</sup> Both components are functions of the flow rate and valve port geometry. Thus, coupling exists between the motions of the spool and actuator load. However, in practical design, the actuating force on the valve spool can be made large compared with  $F_R$  by adequate choice of torque-motor. Furthermore, proper valve design permits reduction of Coulomb and stiction forces to some nominal values. Since the spool motion is small, equation 2 can be approximated by a second-order linear differential equation of the form:

$$\ddot{x}_s(\tau) + 2\zeta_0 \omega_{n0} \dot{x}_s(\tau) + \omega_{n0}^2 x_s(\tau) = \frac{F_T[i_M(\tau)]}{M_s} \quad (3)$$

The damping ratio  $\zeta_0$  and natural frequency  $\omega_{n0}$  may be estimated from the manufacturers' valve frequency response data.

The flow through the valve ports is derived by applying equation 1 and the continuity relationships; see Fig. 1:

$$Q_1 = Q_{s1} - Q_{1r} \quad (4)$$

$$Q_2 = Q_{s2} - Q_{2r} \quad (5)$$

By considering the mass flow rate into the actuator cylinder volumes  $V_1$  and  $V_2$ , relationships between  $Q_1$ ,  $Q_2$ , and the load motion are established.

$$Q_1 = \frac{V_1(\tau)}{\beta} \dot{P}_1(\tau) + \dot{V}_1(\tau) + Q_L \quad (6)$$

Fig. 2 (left). Discontinuous boundaries in pressure ( $P_1$ ,  $P_2$ ) plane

Fig. 3 (right). Normalized velocity response for a relaxed, frictionless system with step-current input

$$Q_2 = \frac{V_2(\tau)}{\beta} \dot{P}_2(\tau) + \dot{V}_2(\tau) - Q_L \quad (7)$$

where  $\beta$  is the bulk modulus of oil. Since the actuator volumes  $V_1(\tau)$  and  $V_2(\tau)$  vary linearly with load displacement—that is  $V_1(\tau) = [V_T/2 + A_a x_L(\tau)]$ ,  $V_2(\tau) = [V_T/2 - A_a x_L(\tau)]$ —equations 6 and 7 reduce to:

$$Q_1 = \frac{\left[ \frac{V_T}{2} + A_a x_L(\tau) \right]}{\beta} \dot{P}_1(\tau) + \dot{x}_L(\tau) A_a + Q_L \quad (8)$$

$$Q_2 = \frac{\left[ \frac{V_T}{2} - A_a x_L(\tau) \right]}{\beta} \dot{P}_2(\tau) - \dot{x}_L(\tau) A_a - Q_L \quad (9)$$

where  $V_T$  is the total oil volume inside the actuator cylinder.

Also, the leakage flow rate  $Q_L$  across the actuator piston can be approximated by a steady flow in an annulus between a finite length circular shaft and a concentric cylinder. In this case,  $Q_L$  is proportional to the pressure differential across the piston.

$$Q_L = C_L [P_1(\tau) - P_2(\tau)] \quad (10)$$

Finally, the equation for load motion is

$$M_L \ddot{x}_L(\tau) = A_a [P_1(\tau) - P_2(\tau)] - F_L[\dot{x}_L(\tau)] + F_D(\tau) \quad (11)$$

$F_D(\tau)$  is an external load disturbance force.

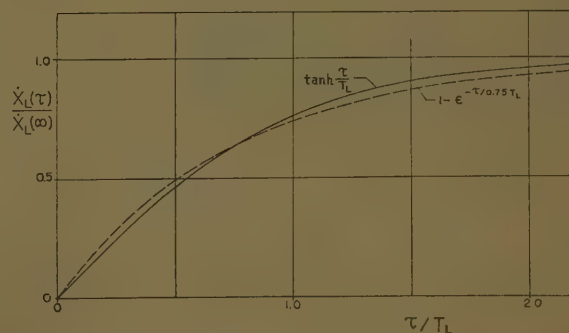
The complete system equations can be conveniently written in the following vector form:

$$\frac{d\mathbf{U}}{d\tau} = \mathbf{H}(\mathbf{U}) + \mathbf{\Psi}(\tau) \quad (12)$$

where  $\mathbf{U}(\tau)$  is the system state vector given by

$$\mathbf{U}(\tau) = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \\ u_6 \end{bmatrix} = \begin{bmatrix} x_s(\tau) \\ \dot{x}_s(\tau) \\ x_L(\tau) \\ \dot{x}_L(\tau) \\ P_1(\tau) \\ P_2(\tau) \end{bmatrix} \quad (13)$$

$\mathbf{\Psi}(\tau)$  is a vector forcing function



$$\Psi(\tau) = \begin{bmatrix} 0 \\ F_T[i_M(\tau)]/M_s \\ 0 \\ F_D(\tau)/M_L \\ 0 \\ 0 \end{bmatrix} \quad (14)$$

and  $\mathbf{H}(\mathbf{U})$  is a vector whose components  $h_1$  are

$$h_1 = \dot{x}_s(\tau) \quad (15)$$

$$h_2 = -2\zeta\omega_n\dot{x}_s(\tau) - \omega_n^2 x_s(\tau) \quad (16)$$

$$h_3 = \dot{x}_L(\tau) \quad (17)$$

$$h_4 = \frac{A_a}{M_L} [P_1(\tau) - P_2(\tau)] - \frac{F_L[\dot{x}_L(\tau)]}{M_L} \quad (18)$$

For "zero-lapped" valve

$$x_s(\tau) > 0 \begin{cases} h_5^+ = \frac{\beta}{\left[\frac{V_T}{2} + A_a x_L(\tau)\right]} \{C_v |A_p[x_s(\tau)]| \operatorname{sgn}[P_s - P_1(\tau)] \sqrt{P_s - P_1(\tau)} - A_a \dot{x}_L(\tau) + C_L[P_1(\tau) - P_2(\tau)]\} \\ h_6^+ = \frac{\beta}{\left[\frac{V_T}{2} - A_a x_L(\tau)\right]} \{-C_v |A_p[x_s(\tau)]| \operatorname{sgn}[P_2(\tau) - P_r] \sqrt{P_2(\tau) - P_r} + A_a \dot{x}_L(\tau) - C_L[P_1(\tau) - P_2(\tau)]\} \end{cases} \quad (19)$$

$$x_s(\tau) < 0 \begin{cases} h_5^- = \frac{\beta}{\left[\frac{V_T}{2} + A_a x_L(\tau)\right]} \{-C_v |A_p[x_s(\tau)]| \operatorname{sgn}[P_1(\tau) - P_r] \sqrt{P_1(\tau) - P_r} - A_a \dot{x}_L(\tau) + C_L[P_1(\tau) - P_2(\tau)]\} \\ h_6^- = \frac{\beta}{\left[\frac{V_T}{2} - A_a x_L(\tau)\right]} \{C_v |A_p[x_s(\tau)]| \operatorname{sgn}[P_s - P_2(\tau)] \sqrt{P_s - P_2(\tau)} + A_a \dot{x}_L(\tau) - C_L[P_1(\tau) - P_2(\tau)]\} \end{cases} \quad (22)$$

Hereafter,  $(\ )^+$  and  $(\ )^-$  notations will be used to indicate the polarity of  $x_s(\tau)$ . Equations for  $h_5$  and  $h_6$  for underlapped and overlapped valves are given in Appendix I.

The system trajectories are restricted to a partially closed region in phase space ( $\mathbf{U}$  space) as a result of the following constraints induced by physical limitations and design specifications:

$$\begin{cases} |x_s(\tau)| \leq x_{s(\max)}: & \text{The spool travel is limited by the valve's physical size.} \\ |\dot{x}_s(\tau)| \leq \dot{x}_{s(\max)}: & \text{The spool velocity is limited by friction and motor saturation.} \\ |x_L(\tau)| \leq x_{L(\max)}: & \text{The actuator piston travel is confined to the working length of the cylinder such that } x_{L(\max)} < V_T/2A_a. \\ |\dot{x}_L(\tau)| \leq \dot{x}_{L(\max)}: & \text{For an external disturbance-free system, the peak actuator piston velocity is limited primarily by friction and supply pressure } P_s. \\ P_{\max} > P_1 \geq P_v \text{ and } P_{\max} > P_2 \geq P_v \} & \text{The instantaneous pressures should be kept below a specified safe pressure } P_{\max}. \text{ Assuming that the oil in the actuator cylinder is ideal and does not completely vaporize, } P_1 \text{ and } P_2 \text{ cannot be lower than the oil vapor pressure } P_v. \end{cases} \quad (23)$$

By inspection of equation 12, it is evident that its right-hand side is discontinuous in nature. Moreover, the phase co-ordinates are bounded due to constraints given by equation 23. The discontinuities arise from the bidirectional flow property of the valve ports and non-symmetry of the valve load.

The system equations may be rewritten in an alternative general vector form:

$$\frac{d\mathbf{U}}{d\tau} = \mathbf{H}_i(\mathbf{U}) + \Psi(\tau), \quad i=1, 2, \dots, N \quad (24)$$

plus a set of mode boundaries

$$Z_k(\mathbf{U}) = 0, \quad k=1, 2, \dots, M \quad (25)$$

Equation 25 represents a set of hypersurfaces which partition the phase space into distinct regions  $\Omega_j$ . Within each  $\Omega_j$ , the system behavior is described by a particular equation in 24, whose  $\mathbf{H}_i(\mathbf{U})$  is continuous in  $\mathbf{U}$  everywhere within  $\Omega_j$ .

For the system under consideration, the mode boundaries are

$$\begin{aligned} x_s(\tau) &= 0 \text{ (for "zero-lapped" valve)} \\ P_s - P_1(\tau) &= 0 \\ P_s - P_2(\tau) &= 0 \\ P_1(\tau) - P_r &= 0 \\ P_2(\tau) - P_r &= 0 \\ \dot{x}_B(\tau) &= 0 \text{ (if Coulomb friction load is present)} \end{aligned} \quad (26)$$

( $P_1, P_2$ ) plane. The scales are greatly exaggerated so as to reveal various mode regions. The permissible trajectories are confined within and to the boundaries of the unshaded regions. Consider a typical pressure phase trajectory, generated by certain forcing  $\Psi(\tau)$ , starting at  $a$ . Its mode sequence is as follows:

Trajectory	Mode Region	Pressure Boundary
ab	I	$P_r < P_1(\tau) \leq P_s$ $P_r < P_2(\tau) \leq P_s$
bc	VII	$P_r < P_1(\tau) \leq P_s$ $P_v < P_2(\tau) \leq P_r$
cd	IV	$P_s \leq P_1(\tau) < P_r$ $P_v \leq P_2(\tau) < P_r$
de	Boundary of IV (cavitation)	$P_s \leq P_1(\tau) < P_r$ $P_2(\tau) = P_v$
ef	Boundary of VII (cavitation)	$P_r < P_1(\tau) \leq P_s$ $P_2(\tau) = P_v$
fa	VII	$P_r < P_1(\tau) < P_s$ $P_v \leq P_2(\tau) \leq P_r$

## Time-Domain Trajectories

In this section, analytical expressions for the time-domain trajectories of a disturbance-free system with "zero-lapped" valve and negligible actuator leakage can be derived. The return pressure  $P_r$  is assumed to be zero since  $P_s \gg P_r$ .

### RESPONSE TO LARGE-STEP MOTOR-CURRENT INPUTS

The large-step motor-current response is of particular interest to design of time-optimal electrohydraulic servomechanisms, where the maximum allowable motor current is applied during the initial mode of operation.

Let the step input current be

$$i_M(\tau) = \begin{cases} 0 & \text{for } \tau < 0 \\ i_0 & \text{for } \tau > 0 \end{cases}$$

Since  $F_T$  is amplitude-sensitive, the effective motor forcing is also a step with magnitude  $F_T(i_0)$ .

The steady-state equations for the system are

$$\begin{aligned} x_s(\infty) &= F_T(i_0)/\omega_n^2 M_s \\ A_a \dot{x}_L(\infty) &= C_v A_p [x_s(\infty)] \sqrt{P_s - P_1(\infty)} \\ A_a \dot{x}_L(\infty) &= C_v A_p [x_s(\infty)] \sqrt{P_2(\infty)} \\ A_a [P_1(\infty) - P_2(\infty)] &= F_L[\dot{x}_L(\infty)] \end{aligned}$$

For viscous friction load,  $F_L[\dot{x}_L(\infty)] = f_L \dot{x}_L(\infty)$ , the steady-state solutions have the form:

$$\begin{aligned} \dot{x}_L(\infty) &= \frac{C_v^2 A_p^2 [x_s(\infty)]}{4A_a^2} \times \left\{ \sqrt{\frac{f_L^2}{A_a^2} + \frac{8A_a^2 P_s}{C_v^2 A_p^2 [x_s(\infty)]}} - \frac{f_L}{A_a} \right\} \\ P_2(\infty) &= \frac{A_a^2 \dot{x}_L^2(\infty)}{C_v^2 A_p^2 [x_s(\infty)]} \\ P_1(\infty) &= P_s - P_2(\infty), \end{aligned}$$

where  $x_s(\infty)$  is given by equation 28.

There are some basic questions pertaining to the behavior of the trajectories upon traversing the mode boundaries. For example, the pressures may change rapidly across certain boundaries. Such cases will not be discussed in the present paper.

From the preceding discussions, it is evident that the electrohydraulic actuator is basically a multiple-mode nonlinear system. This point is further clarified by Fig. 2, showing the projection of a typical phase-space trajectory onto the pressure



Upon initiation of a step current,  $P_1$  and  $P_2$  change rapidly as a result of opening of the valve ports. Since the load cannot move instantaneously, the pressure change is primarily due to compressibility of oil inside the actuator cylinder. Setting  $A_a \ddot{x}_L = 0$  and  $A_a \dot{x}_L(\tau) = A_a \dot{x}_L(0)$  in equations 19 and 20 (for  $i_0 > 0$ ) gives approximate differential equations describing the initial pressure buildup.

$$\approx \left[ \frac{V_T}{2} - A_a x_L(0) \right] \times \{ -C_p A_p [x_s(\tau)] \sqrt{P_2(\tau)} \} \quad (36)$$

$$(\tau) = P_s - \left\{ (P_s/2)^{1/2} + \frac{\beta C_v}{[V_T + 2A_0 x_L(0)]} \times \int_0^\tau A_v [x_s(t)] dt \right\}^2 \quad (37)$$

$$(\tau) = \left\{ (P_s/2)^{1/2} + \frac{\beta C_v}{[V_T - 2A_{ax}L(0)]} \times \int_0^\tau A_p[x_s(t)]dt \right\}^2 \quad (38)$$

$$(\tau) = \frac{F_T(i_0)}{\omega_n v^2 M_s} [1 - e^{-\zeta \omega_n v \tau} \times \cos \omega_n \sqrt{1 - \zeta^2} \tau] \quad (39)$$

After the pressures have reached their peak values, the oil flows are primarily due to actuator piston motion, and the

$$\frac{\left[ \frac{V_T}{2} + A_2 x_L(\tau) \right]}{\beta} \frac{dP_1}{d\tau} \approx 0, \quad \frac{\left[ \frac{V_T}{2} - A_2 x_L(\tau) \right]}{\beta} \frac{dP_2}{d\tau} \approx 0 \quad (40)$$
$$\frac{d\dot{x}_L}{d\tau} = \frac{A_a}{M_L} \left\{ P_s - 2 \left[ \frac{A_a \dot{x}_L(\tau)}{C_v A_p [x_s(\tau)]} \right]^2 \right\} - \frac{f_L \dot{x}_L(\tau)}{M_L} \quad (41)$$

The subsequent system trajectories can be accurately described by the solution of equation 41:

$$\dot{x}_L(t+\tau) = \frac{C_v^2 W_0^2 x_s^2(\infty) M_L}{2 A_a^3 T_L} \tanh \left[ \frac{\tau}{T_L} + \tanh^{-1} \left( \frac{2 A_a^3 T_L \dot{x}_L(t)}{C_v^2 W_0^2 x_s^2(\infty) M_L} + \frac{f_L T_L}{2 M_L} \right) \right] - \frac{f_L C_v^2 W_0^2 x_s^2(\infty)}{4 A_a^3} \quad (42)$$

$$\frac{C_v^2 W_0^2 x_s^2(\infty) M_L}{2A_a^3} \ln \left[ \frac{\cosh \left[ \frac{\tau}{T_L} + \tanh^{-1} \left( \frac{2A_a^2 \dot{T}_L \dot{x}_L(t)}{C_v^2 W_0^2 x_s^2(\infty) M_L} + \frac{f_L T_L}{2M_L} \right) \right]}{\cosh \left[ \tanh^{-1} \left( \frac{2A_a^3 T_L \dot{x}_L(t)}{C_v^2 W_0^2 x_s^2(\infty) M_L} + \frac{f_L T_L}{2M_L} \right) \right]} \right] - \frac{f_L}{2M_L} \tau \quad (43)$$

$$P_1(t+\tau) = P_s - \frac{A_a^2}{C_p^2 W_a^2 x_s^2(\infty)} \dot{x}_L^2(t+\tau) \quad (44)$$

$$P_2(t+\tau) = P_s - P_1(t+\tau) \quad (45)$$

$$T_L = \frac{C_v W_0 x_s(\infty) M_L}{2A a^2} \times \left[ \frac{P_s}{2} + \frac{f_L^2 C_v^2 W_0^2 x_s^2(\infty)}{16 a^4} \right]^{-1/2} \quad (46)$$

A normalized plot of load velocity given by equation 42 with  $\dot{x}_L(t)$  and  $f_L$  equal to zero is shown in Fig. 3. For this case, the step response of a relaxed, frictionless system for a particular current  $i_0$  can be approximated by that of a one-integral plus one-time-constant system. The equivalent time constant  $T'_L$  is proportional to the input current amplitude  $i_0$  and has a value of approximately  $0.75 T_L$ . Equations 42 and 43 are useful

Power series extrapolations of system trajectories in the time domain have been used in the design of predictor control systems with linear dynamic processes.<sup>12</sup> This approach is especially useful in determining the required forcing functions if the control system inputs can be also extrapolated by power series in time. For nonlinear processes such as the electrohydraulic actuator considered in this paper, the use of a linearized model for accurate prediction of the time-domain trajectories is usually unsatisfactory. The power series approach is applicable.

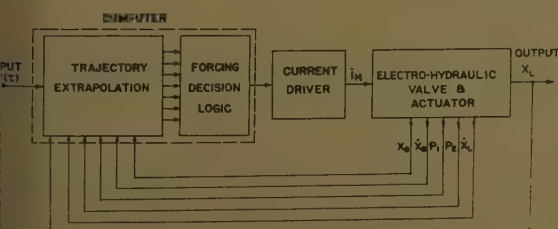
Fig. 4 shows a block diagram of a proposed predictor electrohydraulic servomechanism. The system positional input  $r(\tau)$  and the state variables  $u_i$  of the hydraulic actuator are measured periodically and fed into the computer as initial conditions for a particular computation period. The computer performs extrapolations of both the input,  $r(\tau)$ , and the actuator trajectories corresponding to certain assumed forcings for a future time  $\tau$ . Then, on the basis of the extrapolated information and a predetermined forcing criterion, a control signal is generated and fed into the current driver.

In the subsequent discussion, power series solutions of equation 12 will be considered. Since the valve-spool motion is assumed to be independent of the remaining system dynamics, equation 12 can be written as a time-dependent equation by specifying the port area function as a power series in  $\tau$ .

$$\frac{dY}{d\tau} = G(Y, \tau) \quad (47)$$

where

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} x_L(\tau) \\ \dot{x}_L(\tau) \\ P_1(\tau) \\ P_2(\tau) \end{bmatrix}$$



**Fig. 4 (left). Block diagram of a proposed predictor electro-hydraulic servomechanism**

$$\mathbf{G}(\mathbf{Y}, \tau) = \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \end{bmatrix} = \begin{bmatrix} h_3 \\ h_4 \\ h_5^\pm \\ h_6^\pm \end{bmatrix} \quad (48)$$

$$\text{with } A_p[x_s(\tau)] = A_p(t) + \sum_{n=1}^{\infty} \alpha_n \tau^n \quad (49)$$

To determine power series solutions for equation 47, the system's operating mode region,  $\Omega_j$ , in phase space, must be determined from the initial state  $\mathbf{Y}(t)$  so that the appropriate form of  $\mathbf{G}(\mathbf{Y}, \tau)$  is used for computation. Within  $\Omega_j$ , the components of  $\mathbf{G}$  (that is,  $g_i(\mathbf{Y}, \tau)$ ) are analytic functions of  $x_L, \dot{x}_L, P_1$  and  $P_2$  simultaneously. The general forms of power series solutions are

$$\left. \begin{aligned} x_L(t+\tau) &= x_L(t) + \sum_{n=1}^{\infty} \gamma_{1n} \tau^n \\ \dot{x}_L(t+\tau) &= \dot{x}_L(t) + \sum_{n=1}^{\infty} \gamma_{2n} \tau^n \\ P_1(t+\tau) &= P_1(t) + \sum_{n=1}^{\infty} \gamma_{3n} \tau^n \\ P_2(t+\tau) &= P_2(t) + \sum_{n=1}^{\infty} \gamma_{4n} \tau^n \end{aligned} \right\} \quad (50)$$

where

$$\begin{aligned} \gamma_{11} &= g_1(t) \\ \gamma_{12} &= 1/2 \left[ \sum_{j=1}^4 \frac{\partial g_1}{\partial y_j} \gamma_{j1} + \frac{\partial g_1}{\partial \tau} \right] \mathbf{Y}(\tau=0), \tau=0 \\ \gamma_{13} &= 1/3 \left[ \sum_{j=1}^4 \frac{\partial g_1}{\partial y_j} \gamma_{j2} + 1/2 \sum_{j=1}^4 \times \right. \\ &\quad \left. \sum_{k=1}^4 \frac{\partial^2 g_1}{\partial y_j \partial y_k} \gamma_{j1} \gamma_{k1} + \sum_{j=1}^4 \times \right. \\ &\quad \left. \frac{\partial^2 g_1}{\partial y_j \partial \tau} \gamma_{j1} + 1/2 \frac{\partial^2 g_1}{\partial \tau^2} \right] \mathbf{Y}(\tau=0), \tau=0 \\ &\vdots \\ \gamma_{1N} &= \frac{1}{N} P_{N1}[\gamma_{j1}, \gamma_{j2}, \gamma_{j3}, \dots, \gamma_{j(N-1)}] \\ &\quad j=1 \text{ to } 4 \quad (51) \end{aligned}$$

$P_{N1}$  are polynomials in  $\gamma_{j1}, \gamma_{j2}, \dots, \gamma_{j(N-1)}$ .

The domain of convergence may be estimated by the following inequality:

$$|\tau| < T'(1 - e^{-\mu/6ST'}) \quad (52)$$

The constants  $S, \mu$ , and  $T'$  are chosen so that  $S \geq |g_i|, \mu > \mu_i$ , and  $T > T'$  for all values of  $y_i$  and  $\tau$ , in a closed domain, defined by

$$|y_i(t+\tau) - y_i(t)| \leq \mu_i; |\tau| \leq T, i=1, 2, \dots, 4 \quad (53)$$

where  $\mu_i > 0, T > 0$ , and  $y_i$ 's belong to a particular  $\Omega_j$ , which implies that the system trajectories remain in the same region,  $\Omega_j$ , in phase space during time interval  $[t, t+\tau]$ .

For short extrapolation time  $\tau$ , equation 50 may be simplified by series truncation and by considering the oil volumes on each side of the actuator piston as constants during the extrapolation time interval  $[t, t+\tau]$ .

$$\left. \begin{aligned} V_1(\tau) &= [V_T/2 + A_a x_L(\tau)] \approx \\ &\quad [V_T/2 + A_a x_L(t)] \\ V_2(\tau) &= [V_T/2 - A_a x_L(\tau)] \approx \\ &\quad [V_T/2 - A_a x_L(t)] \end{aligned} \right\} \quad (54)$$

An estimate of truncation error is obtainable from the following inequality for the upper bound of the remainder  $R_k$  for a  $k$ -term truncated series solution:

$$|R_k| \leq \frac{\mu |\tau|^{k+1}}{(T')^k (T' - |\tau|)} \times \left\{ 2 + \left[ \frac{5MT'}{\mu} \ln \left( 1 - \frac{T'}{T'} \right) \right]^{1/6} \right\} \quad (55)$$

$$\text{where } |\tau| < T' < T'(1 - e^{-\mu/6ST'}) \quad (56)$$

It should be remarked that many well-known numerical procedures for starting the solution of a system of differential equations (e.g., Runge-Kutta method) consist of matching the first few terms of the power series solutions. Moreover, in these procedures, no attempt is made to obtain solutions explicit in the independent variable  $\tau$ , but rather to compute the change in the solution due to a given

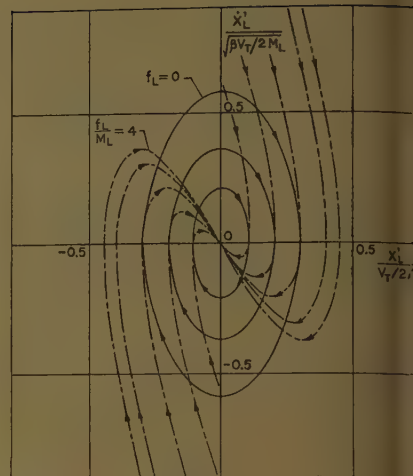


Fig. 5. Phase-plane trajectories of actual load after termination of a rectangular current pulse  $[P_1(T_p) = P_2(T_p)]$

increment in  $\tau$ . The power series solutions, explicitly expressed in terms of initial state variables,  $y_i(t)$ , and extrapolation time  $\tau$ , are useful in determining required forcing function to satisfy specified performance criterion. For example, let the desired load position,  $r(t+\tau)$ , and velocity,  $\dot{r}(t+\tau)$ , within a finite time interval,  $(t, t+\tau)$ , be specified in truncated power series form:

$$r(t+\tau) = r(t) + \dot{r}(t)\tau + \sum_{n=2}^N \Gamma_n(t) \tau^n$$

$$\dot{r}(t+\tau) = \dot{r}(t) + \sum_{n=2}^N n \Gamma_n(t) \tau^{n-1} \quad (57)$$

Then, a possible approach to the control problem is to find the required port coefficients,  $\alpha_n$ , in equation 49 so that prescribed performance criterion defined over the finite time interval,  $[t, t+\tau]$ , may be satisfied. A typical class of performance criteria has the form:

$$\text{minimize } \int_0^\tau \phi[r(t+\tau') - x_L(t+\tau'), \dot{r}(t+\tau') - \dot{x}_L(t+\tau')] d\tau' \quad (58)$$

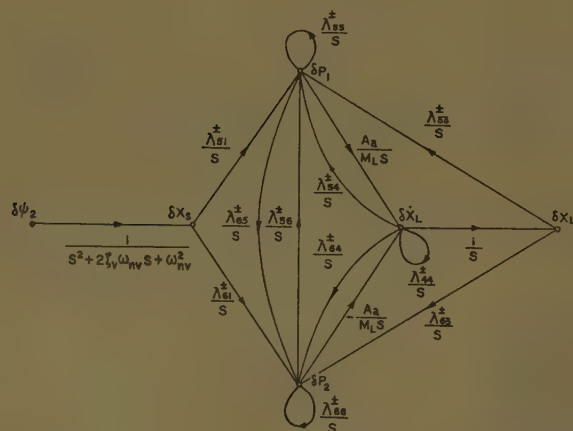
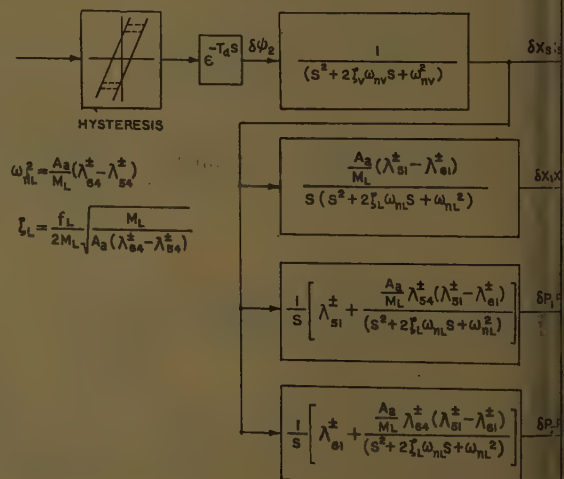


Fig. 6 (left). Flow graph for linearized system

Fig. 7 (right). Transfer function block diagram for a simplified system





where  $\phi$  is a specified function of its arguments.

## RESPONSE TO A RECTANGULAR CURRENT PULSE

In time-optimal and pulse-width-modulated sampled-data control systems, the required forcing of the output dynamic member often takes the form of a time sequence of positive and negative rectangular pulses.<sup>13</sup> The valve-actuator response to such a current-pulse sequence can be readily deduced from the subsequent discussion on single-pulse response.

First, consider the case where the valve-spool motion is instantaneous with respect to torque-motor current. The valve ports open upon pulse initiation at time,  $\tau = 0$ , and close upon pulse termination at  $\tau = T_p$ . Hence the system response during the pulse duration,  $0 < \tau < T_p$ , is identical to the step response, only the trajectories after port closure remain to be determined. Clearly, they correspond to the free-oscillation trajectories of actuator load and the closed oil columns between the actuator piston and the closed valve ports. Setting  $A_p[x_s(\tau)]$ , in equations 19 and 20, equal to zero and performing necessary integrations, the equation for load motion becomes

$$\ddot{x}_L'(\tau') + \frac{F_L[\dot{x}_L'(\tau')]}{M_L} + \frac{A_a\beta}{M_L} \ln \left\{ \left[ 1 + \frac{x_L'(\tau')}{\frac{V_T}{2A_a} + x_L(T_p)} \right] \left[ 1 - \frac{x_L'(\tau')}{\frac{V_T}{2A_a} - x_L(T_p)} \right] \right\} = \frac{A_a}{M_L} [P_1(T_p) - P_2(T_p)] \quad (59)$$

where  $\tau' = \tau - T_p$  and  $x_L'(\tau' = 0) = x_L(T_p)$ . The closed oil columns are in effect non-near hard springs.

Equation 59 is valid if the pressures  $P_1$  and  $P_2$  never drop down to the oil vapor pressure  $P_v$ .

The first integral for equation 59 with  $\dot{x}_L = 0$  is

$$\begin{aligned} & \{[\dot{x}_L'(\tau')]^2 - [\dot{x}_L(T_p)]^2\} \\ &= \frac{A_a}{M_L} \{x_L'(\tau') [P_1(T_p) - P_2(T_p)]\} - \\ & \frac{A_a\beta}{M_L} \ln \left\{ \left( 1 + \frac{x_L'(\tau')}{\left[ \frac{V_T}{2A_a} + x_L(T_p) \right]} \right)^{\left[ \frac{V_T}{2A_a} + x_L(T_p) + x_L'(\tau') \right]} \right. \\ & \quad \times \left. \left( 1 - \frac{x_L'(\tau')}{\left[ \frac{V_T}{2A_a} - x_L(T_p) \right]} \right)^{\left[ \frac{V_T}{2A_a} - x_L(T_p) - x_L'(\tau') \right]} \right\} \quad (60) \end{aligned}$$

For  $F_L[\dot{x}_L'(\tau')] \neq 0$ , the trajectories may be constructed in the phase plane by using the isocline method. Typical trajectories for damped and undamped systems with  $P_1(T_p) = P_2(T_p)$  are shown in Fig. 5. In practical systems, the single-pulse test is useful for experimental deter-

mination of total actuator load damping provided that the valve is quick-acting.

A more realistic case would be to describe the port areas closure by a truncated power series in  $\tau'$ .

$$A_p[x_s(\tau')] = A_p(T_p) + \sum_{n=1}^N \alpha_n(\tau')^n \quad (61)$$

For a value with constant port width,  $W_o$ , and the spool motion describable by equation 3 with  $F_T = 0$ , (since  $\dot{i}_M = 0$  after pulse termination), equation 61 takes the form:

$$A_p[x_s(\tau')] = W_o \left\{ x_s(T_p) + \dot{x}_s(T_p)\tau' - [2\xi\omega_n\tau\dot{x}_s(T_p) + \omega_n^2x_s(T_p)]\frac{\tau'^2}{2} + \dots \right\} \quad (62)$$

Approximate expressions for load trajectories are obtainable by direct application of equations 50 and 51.

$$\dot{x}_L(\tau') = \dot{x}_L(T_p) + \sum_{n=1}^{\infty} \gamma_{2n}(\tau')^n \quad (63)$$

$$x_L(\tau') = x_L(T_p) + \dot{x}_L(T_p)\tau' + \sum_{n=1}^{\infty} \frac{\gamma_{2n}}{n} (\tau')^{n+1} \quad (64)$$

Here,  $\tau'$  must be restricted to the non-cavitating range.

For the case where  $F_L = 0$ ,  $\dot{x}_L(T_p) = C_o A_p(T_p) \sqrt{P_s/2}/A_a$ , and  $P_1(T_p) = P_2(T_p) = P_s/2$ , the first few series coefficients are

$$\begin{aligned} \gamma_{21} &= 0 \\ \gamma_{22} &= 0 \\ \gamma_{23} &= \frac{A_a\beta C_o \alpha_1}{6M_L} \sqrt{\frac{P_s}{2}} \left[ \frac{1}{V_1(T_p)} + \frac{1}{V_2(T_p)} \right] \end{aligned}$$

$$\gamma_{24} = \frac{A_a}{12M_L} \left\{ \beta \alpha_2 C_o \sqrt{\frac{P_s}{2}} \left[ \frac{1}{V_1(T_p)} + \frac{1}{V_2(T_p)} \right] - \frac{\beta^2 C_o^2 \alpha_1 A_p(T_p)}{4} \times \left[ \frac{1}{V_1^2(T_p)} + \frac{1}{V_2^2(T_p)} \right] \right\} \quad (65)$$

It can be seen that if  $\alpha_1 = W_o \dot{x}_s(T_p) = 0$ , then  $\gamma_{23}$  is identically zero. Hence the series solution for  $\dot{x}_L(\tau')$  starts with  $\dot{x}_L(T_p)$  and a term of the order  $(\tau')^4$ . The case just cited is a rather special one. In other situations, the first three coefficients never vanish simultaneously.

The results derived in this section have the following applications:

1. They may be applied to near time-optimal electrohydraulic servomechanism design. Since the oil column between the closed-valve ports and the actuator piston provides an effective means of braking the load motion, a possible operating mode for achieving a near time-optimal noncavitating response (for step positioning) is to open the valve ports as wide as possible during the initial acceleration period and to follow this by a controlled closure of the valve ports. The valve port itself acts as a nonlinear damper. The approximate port closure time function may assume the form of equation 59. By specifying the desired load response during port closure in the form of a polynomial in time, the required coefficients  $\alpha_n$  may be computed.

2. Calculation of system response to a time sequence of rectangular pulses may

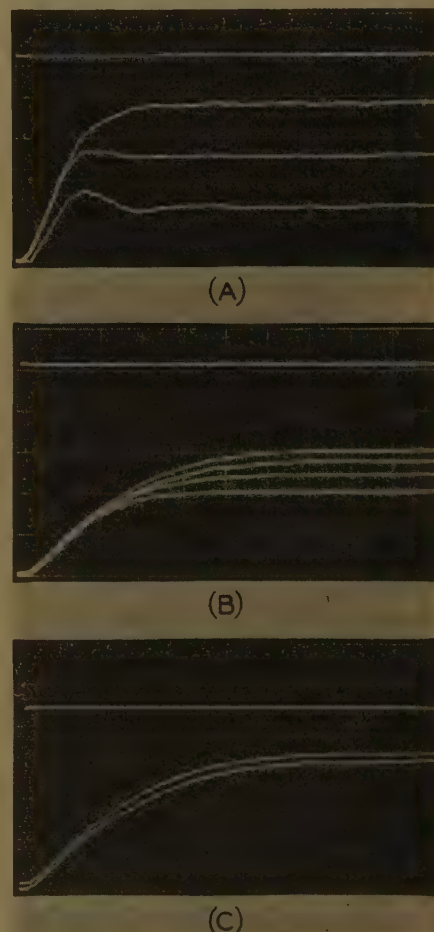


Fig. 8. Velocity response of the physical model to step-current inputs with a time scale of 5 milliseconds per large division and a velocity scale of 4.64 (A) and 11.6 inches per second per large division (B, C), current increments: 1 mA

be pursued in a piece-wise manner using the results for single-pulse response.

### Model for Small Motions

Linearization techniques are applicable for describing the small motion of non-linear systems and investigating their local stability about their equilibrium states. Linearized models of valve-actuator systems have been discussed extensively in the past.<sup>1-3,10</sup> Here, a generalized linear model is derived systematically from the basic system equation given by 12. Due to discontinuities in the right-hand side of equation 12, the usual linearization procedure must be applied with caution, particularly in the case where the system's equilibrium state lies on one of the mode boundaries given by equation 25. In a precision servomechanism with a "zero-lapped" valve and step position inputs, a mathematical model linearized about  $x_s(\tau)=0$  is useful for small-motion and stability studies of the closed-loop system. In the case where the input to the servomechanism is of ramp nature, the torque-motor operates with a bias current,  $i_{Mo}$ , and the steady-state load motion will have a velocity component. Then the linearized model is derived by considering the load position,  $x_L(\tau)$ , quasi-stationary in time.

Let the equilibrium and perturbed state vectors, denoted by  $U_e^\pm$  and  $\delta U^\pm$  respectively, be defined as

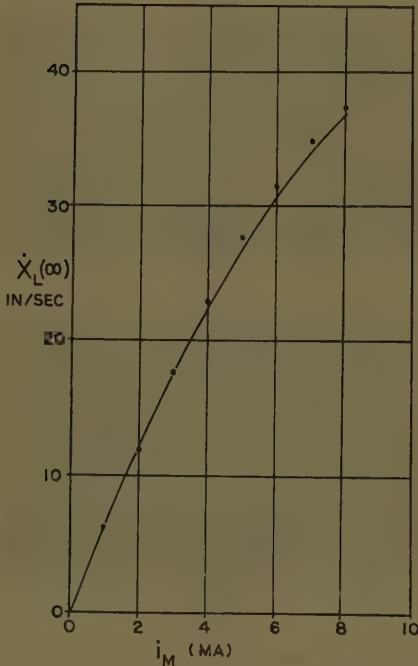


Fig. 9. Experimental and calculated steady-state velocities of the physical model with step-current inputs ( $f_L=1.6$  pounds per inch per second)

$$U_e^\pm = \begin{bmatrix} x_{se}^\pm \\ \dot{x}_{se}^\pm \\ x_{Le}^\pm \\ \dot{x}_{Le}^\pm \\ P_1^\pm \\ P_2^\pm \end{bmatrix}, \delta U = U - U_e^\pm = \begin{bmatrix} \delta x_s \\ \delta \dot{x}_s \\ \delta x_L \\ \delta \dot{x}_L \\ \delta P_1 \\ \delta P_2 \end{bmatrix} \quad (66)$$

The  $( )^+$  and  $( )^-$  notations are used again to denote the system state corresponding to positive and negative  $x_{se}$  respectively. For  $x_{se}=0$ , the notations indicate the sign of  $\delta x_s$ .

Linearization of equation 12 about the equilibrium state leads to

$$\frac{d(\delta U^+)}{d\tau} = \Lambda^+ \delta U^+ + \delta \Psi^+(\tau) \quad (67)$$

and

$$\frac{d(\delta U^-)}{d\tau} = \Lambda^- \delta U^- + \delta \Psi^-(\tau) \quad (68)$$

where  $\Lambda^+$  and  $\Lambda^-$  are the Jacobian matrices, associated with  $H(U)$ , evaluated at  $U_e^+$  and  $U_e^-$  respectively. Their explicit forms are

$$\Lambda^\pm = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -\omega_{nv}^2 & -2\zeta_v \omega_{nv} & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & \lambda_{44}^\pm \\ \lambda_{51}^\pm & 0 & \lambda_{53}^\pm & \lambda_{54}^\pm \\ \lambda_{51}^\pm & 0 & \lambda_{53}^\pm & \lambda_{54}^\pm \end{bmatrix}$$

$$\lambda_{ij}^+ = \frac{\partial h_i^+}{\partial u_j^+} \Big|_{U_e^+}, \lambda_{ij}^- = \frac{\partial h_i^-}{\partial u_j^-} \Big|_{U_e^-}$$

$\delta \Psi^+(\tau)$  is a perturbed forcing vector by

$$\delta \Psi^+(\tau) = \begin{bmatrix} 0 \\ \delta \psi_2(\tau) \\ 0 \\ \delta \psi_4(\tau) \\ 0 \\ 0 \end{bmatrix} = \Psi^+(\tau) - \begin{bmatrix} 0 \\ F_T(i_{Mo})/M_s \\ 0 \\ F_{Do}/M_L \\ 0 \\ 0 \end{bmatrix} \quad (70)$$

where  $F_{Do}$  is the mean value of the external load disturbance force.

Performing a Laplace transformation of equations 67 and 68 leads to

$$\begin{aligned} (sI - \Lambda^+) \delta U^+(s) &= \delta \Psi^+(s) \\ (sI - \Lambda^-) \delta U^-(s) &= \delta \Psi^-(s) \end{aligned} \quad (71)$$

Equivalent transfer functions relating  $\delta \psi_2(s)$  and  $\delta \psi_4(s)$  with various components of  $\delta U^\pm$  are obtainable by reduction of the flow graph of Fig. 6 associated with equation 71.

For a symmetrical, "zero-lapped," rectangular-port (width  $W_o$ ) valve with viscous load friction and negligible leakage, the equilibrium state vector for  $i_{Mo}=0, (x_{se}=0)$  is

$$U_e = \begin{bmatrix} 0 \\ 0 \\ x_{Le} \\ 0 \\ P_s/2 \\ P_s/2 \end{bmatrix} \quad (72)$$

The matrix elements  $\lambda_{ij}^\pm$  are

$$\begin{aligned} \lambda_{44}^\pm &= -\frac{f_L}{M_L} \\ \lambda_{51}^\pm &= \frac{\beta C_v W_o}{\left(\frac{V_T}{2} + A_a x_{Le}\right)} \sqrt{\frac{P_s}{2}} \\ \lambda_{61}^\pm &= \frac{-\beta C_v W_o}{\left(\frac{V_T}{2} - A_a x_{Le}\right)} \sqrt{\frac{P_s}{2}} \\ \lambda_{54}^\pm &= \frac{-\beta A_a}{\left(\frac{V_T}{2} + A_a x_{Le}\right)} \\ \lambda_{64}^\pm &= \frac{\beta A_a}{\left(\frac{V_T}{2} - A_a x_{Le}\right)} \\ \lambda_{55}^\pm &= \lambda_{66}^\pm = \lambda_{55}^\pm = \lambda_{66}^\pm \\ &= \lambda_{33}^\pm = \lambda_{63}^\pm = 0 \end{aligned}$$

As a result of valve symmetry, the corresponding matrix elements for  $\delta x_s > 0$

$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ A_a & -A_a \\ M_L & M_L \\ \lambda_{55}^\pm & \lambda_{56}^\pm \\ \lambda_{65}^\pm & \lambda_{66}^\pm \end{bmatrix} \quad (6)$$

and  $\delta x_s < 0$  are identical, so that the system's small motion about  $x_s=0$  can be described by a single vector equation. For a valve with nonsymmetrical flow characteristics, two equations in the forms given by 71 are necessary for a complete description.

A transfer function diagram is shown in Fig. 7. It can be seen from the results that the damping of the actuator pole is governed by load friction force alone. This is expected, since for small oil flow the damping effects of the valve ports are negligible. Valve underlap and leakage across the actuator piston generally introduce additional damping. Hysteresis and dead-time in the valve torque-motor are included for stability analysis. Hence the well-known describing function techniques may be used.

### Experimental Studies

To check the accuracy of the derived mathematical models, experimental tests were performed on a typical valve-controlled actuator system and its corresponding mathematical model simulated by an analog computer. Pertinent system parameters are listed in Appendix II.

First, the step-current response of the system was determined. Fig. 8 shows



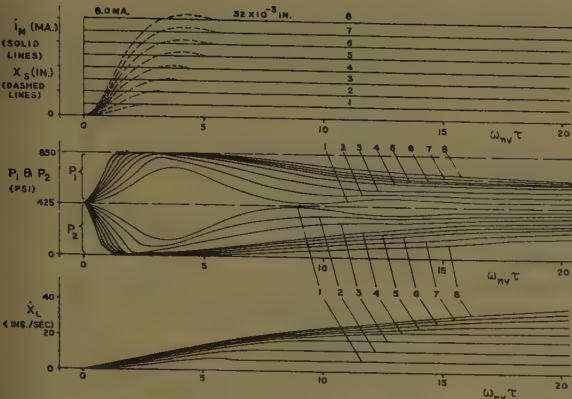


Fig. 10. Analog computer solutions of system's step-current response ( $\omega_{nv}=943$  radians per second)

the velocity responses for various step-current amplitudes. Fig. 9 shows that the steady-state velocities are in close agreement with the values computed using equation 32. The viscous friction coefficient  $f_L$  used in the above calculations was experimentally estimated from single rectangular pulse test results. The velocity profiles for large currents agree within  $\pm 5\%$  with the corresponding analog computer solutions in Fig. 10 and the approxi-

mate analytical result given by equation 42. For small current amplitudes, less than 2.0 ma (milliamperes), the velocity responses of the physical model become less damped as a result of the weaker damping effect of the valve ports, whereas this effect is predominant in large current response.

Figs. 11 and 12 show the system responses for single- and multiple-pulse inputs. Actuator pressures  $P_1$  and  $P_2$  are monitored by means of pressure transducers. The multiple mode nature of the

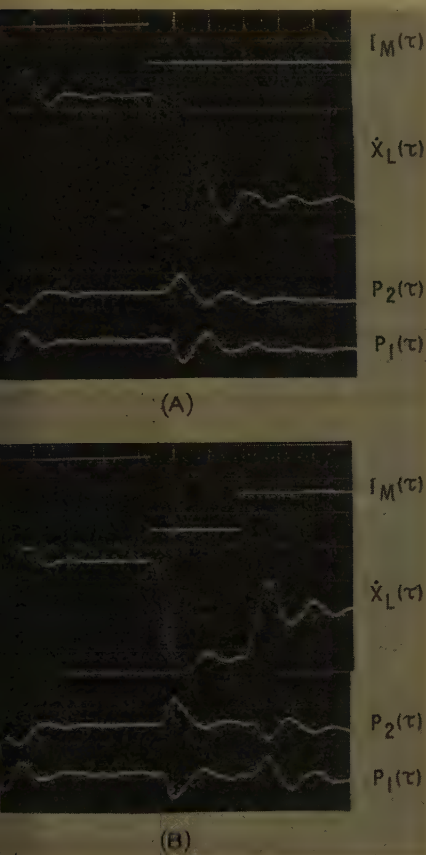


Fig. 11. (A) Single current-pulse response of the physical model with a pulse amplitude of 1 ma and a velocity scale of 2.32 inches per sec per large division. (B) Double current-pulse response of the physical model with a pulse amplitude of +1 ma and a velocity scale of 4.64 inches per second per large division. Time scale: 10 milliseconds per large division, pressure scale: 708 pounds per square inch per large division

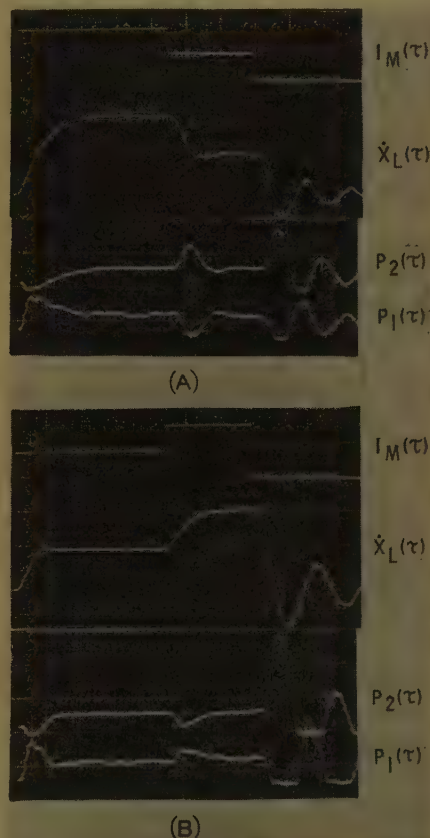


Fig. 12. Response of the physical model to a sequence of current pulses with a velocity scale of 11.6 inches per second per large division, a pressure scale of 708 pounds per square inch per large division, and pulse sequences of +4, +2, 0 ma (A) and +2, +4, 0 ma (B)

system is clearly revealed in the lower portion of Fig. 12. The difference in rise times is shown by the velocity responses to step currents, +2 and +4 ma, which agree with the approximate results predicted by equation 42. After the current has returned to zero, the oscillatory response is due to the motion of the actuator load coupled with the closed oil columns between the valve port and actuator piston. High-pressure peaks exceeding  $P_s$  and cavitation are clearly indicated.

## Conclusions

It has been shown that an electro-hydraulic valve controlled actuator can be represented by reasonably simple and accurate mathematical models by making suitable assumptions. The basic model is shown to be a multiple-mode nonlinear system. Simple results are derived for step and pulsed valve motor-current inputs and their applications are discussed. Finally, the adequacy of the models for time-domain design is substantiated by experimental studies of a physical model and its corresponding mathematical model simulated by an analog computer.

## Appendix I. $h_5$ and $h_6$ for Underlapped and Overlapped Valves

### Underlapped Valve (With Lapping $\pm \Delta_u$ )

For  $|x_s(\tau)| < \Delta_u$

$$h_5 = \left[ \frac{\beta}{\frac{V_T}{2} + A_a x_L(\tau)} \right] \{ C_v' |A_p x_s(\tau)| \times [\text{sgn}[P_s - P_1(\tau)]\sqrt{|P_s - P_1(\tau)|} - \text{sgn}[P_1(\tau) - P_r]\sqrt{|P_1(\tau) - P_r|} - A_a \dot{x}_L(\tau) + C_L[P_1(\tau) - P_2(\tau)]] \} \quad (74)$$

$$h_6 = \left[ \frac{\beta}{\frac{V_T}{2} - A_a x_L(\tau)} \right] \{ C_v' |A_p x_s(\tau)| \times [\text{sgn}[P_s - P_2(\tau)]\sqrt{|P_s - P_2(\tau)|} - \text{sgn}[P_2(\tau) - P_r]\sqrt{|P_2(\tau) - P_r|} + A_a \dot{x}_L(\tau) - C_L[P_1(\tau) - P_2(\tau)]] \} \quad (75)$$

$C_v'$  is a discharge proportionality constant for small port opening.

For  $x_s(\tau) > \Delta_u$ : same as equations 19 and 20

For  $x_s(\tau) < -\Delta_u$ : same as equations 21 and 22

### Overlapped Valve (With Lapping $\pm \Delta_o$ )

For  $|x_s(\tau)| < \Delta_o$

$$h_s = \frac{-\beta A_a \dot{x}_L(\tau)}{\left[ \frac{V_T}{2} + A_a x_L(\tau) \right]} \quad (76)$$

$$h_s = \frac{\beta A_a \dot{x}_L(\tau)}{\left[ \frac{V_T}{2} - A_a x_L(\tau) \right]} \quad (77)$$

For  $x_s(\tau) > \Delta_o$

$$h_s = \frac{\beta}{\left[ \frac{V_T}{2} + A_a x_L(\tau) \right]} \{ C_v A_p [x_s(\tau) - \Delta_o] \times \text{sgn} [P_s - P_1(\tau)] \sqrt{|P_s - P_1(\tau)|} - A_a \dot{x}_L(\tau) + C_L [P_1(\tau) - P_2(\tau)] \} \quad (78)$$

$$h_s = \frac{\beta}{\left[ \frac{V_T}{2} - A_a x_L(\tau) \right]} \{ -C_v A_p [x_s(\tau) - \Delta_o] \times \text{sgn} [P_2(\tau) - P_r] \sqrt{|P_2(\tau) - P_r|} + A_a \dot{x}_L(\tau) - C_L [P_1(\tau) - P_2(\tau)] \} \quad (79)$$

For  $x_s(\tau) < -\Delta_o$

$$h_s = \frac{\beta}{\left[ \frac{V_T}{2} + A_a x_L(\tau) \right]} \{ -C_v A_p [x_s(\tau) + \Delta_o] \times \text{sgn} [P_1(\tau) - P_r] \sqrt{|P_1(\tau) - P_r|} - A_a \dot{x}_L(\tau) + C_L [P_1(\tau) - P_2(\tau)] \} \quad (80)$$

$$h_s = \frac{\beta}{\left[ \frac{V_T}{2} - A_a x_L(\tau) \right]} \{ C_v A_p [x_s(\tau) + \Delta_o] \times \text{sgn} [P_s - P_2(\tau)] \sqrt{|P_s - P_2(\tau)|} + A_a \dot{x}_L(\tau) - C_L [P_1(\tau) - P_2(\tau)] \} \quad (81)$$

## Appendix II. Parameters of the Physical Model

### Actuator

Inertial load ( $M_L$ ): 0.0311 pound second squared per inch (12-pound weight)  
Load Coulomb friction force ( $F_c$ ): approximately 4 pounds  
Load viscous friction force coefficient  $f_L$ , determined from single-pulse response: 1.6 pounds per inch per second  
Piston area ( $A_a$ ): 0.125 inch squared  
Maximum piston stroke: 5.0 inches  
Total volume of oil in cylinder ( $V_T$ ): 0.99 inch<sup>3</sup>  
Bulk modulus of oil ( $\beta$ ):  $2.5 \times 10^5$  pounds per square inch

### Valve (2-Stage Flapper Type)

$\omega_{np} = 943$  radians per second (150 cycles per second)  
 $\zeta_v = 0.7$   
Port width ( $W_o$ ): 0.1 inch  
Spool displacement for  $i_M = 4.0$  ma: 0.016 inch  
Port discharge proportionality constant: 104 inches<sup>2</sup>/(pound)<sup>1/2</sup> second

### Hydraulic Supply

Supply pressure ( $P_s$ ): 850 pounds per square inch (gage)  
Return pressure ( $P_r$ ): atmospheric

## References

1. FLUID POWER CONTROL (book), J. F. Blackburn, G. Reethof, J. L. Shearer, editors. John Wiley & Sons, Inc., New York, N. Y., 1960.
2. DYNAMIC OPERATION OF A FORCE-COMPEN-

- SATED HYDRAULIC THROTTLING VALVE, J. L. Bower and F. B. Tuteur. *Transactions, American Society of Mechanical Engineers*, New York, N. Y., vol. 75, 1953, pp. 1395-1406.
3. AN ANALYSIS OF THE DYNAMICS OF HYDRAULIC SERVOMOTORS UNDER INERTIA LOADS AND THE APPLICATION TO DESIGN, H. Gold, E. W. Otto, V. L. Ransom. *Ibid.*, pp. 1383-94.
4. A THEORETICAL ANALYSIS OF THE RESPONSE OF A LOADED HYDRAULIC RELAY, R. Butler. *Proceedings, Institute of Mechanical Engineers*, London, England, vol. 173, no. 16, 1959, pp. 429-58.
5. THE RESPONSE OF A LOADED HYDRAULIC SERVOMECHANISM, D. E. Turabull. *Ibid.*, no. 9, 1959, pp. 270-93.
6. A DESCRIBING FUNCTION FOR MULTIPLE NONLINEARITIES PRESENT IN ELECTROHYDRAULIC CONTROL VALVES, J. Zaborsky, H. J. Harrington. *AIEE Transactions*, pt. I (Communication and Electronics), vol. 76, May 1957, pp. 183-90.
7. GENERALIZED CHARTS OF THE EFFECTS OF NONLINEARITIES IN ELECTROHYDRAULIC CONTROL VALVES, J. Zaborsky, H. J. Harrington. *Ibid.*, pp. 191-98.
8. A DESCRIBING FUNCTION FOR THE MULTIPLE NONLINEARITIES PRESENT IN 2-STAGE ELECTROHYDRAULIC CONTROL VALVES, J. Zaborsky, H. J. Harrington. *AIEE Transactions*, pt. II (Applications and Industry), vol. 76, Jan. 1958, pp. 394-411.
9. GENERALIZED CHARTS OF THE EFFECTS OF NONLINEARITIES IN 2-STAGE ELECTROHYDRAULIC CONTROL VALVES, J. Zaborsky, H. J. Harrington. *Ibid.*, pp. 401-08.
10. THE ANALYSIS OF VALVE-CONTROLLED HYDRAULIC SERVOMECHANISMS, R. G. Rausch. *Brooklyn System Technical Journal*, New York, N. Y., vol. 38, 1959, pp. 1513-60.
11. CONTRIBUTIONS TO HYDRAULIC CONTROL PART I AND II, S. Y. Lee, J. F. Blackburn. *Transactions, American Society of Mechanical Engineers*, vol. 74, 1952, pp. 1005-16.
12. RELAY-TYPE CONTROL SYSTEMS DESIGNED FOR RANDOM INPUTS, A. M. Hopkin, P. K. C. Wang. *AIEE Transactions*, pt. II (Applications and Industry), vol. 78, Sept. 1959, pp. 227-33.
13. PULSE WIDTH CONTROL OF SAMPLED-DATA SYSTEMS, W. L. Nelson. *Technical Report 35/B*, Dept. of Elec. Eng., Columbia University, New York, N. Y., July 1959.

# Signal Stabilization of a Control System with Random Inputs

R. OLDENBURGER  
MEMBER AIEE

R. SRIDHAR  
NONMEMBER AIEE

**N**OISE AND ITS EFFECT on system performance now can be predetermined by a designer if he uses the theory advanced herein. An external signal is introduced to stabilize a nonlinear feedback system in a state of self-sustained oscillation, but the introduction should be made at a sufficiently high frequency and at a convenient point in the loop. This phenomenon, first introduced by R. Oldenburger, one of the authors, was termed "signal stabilization."<sup>1</sup>

This author, in collaboration with Liu,<sup>2</sup> later explained signal stabilization in the case of a system with one nonlinear element when the external or

stabilizing signal was sinusoidal. They defined an equivalent gain for a nonlinearity and showed that the latter can often be replaced by the former for stability analysis. Equivalent gain was described as the "limit of the output divided by the average value of the input

Paper 61-712, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE-AICHE-ASME-IRE-ISA Joint Automatic Control Conference, Boulder, Colo., June 28-30, 1961. Manuscript submitted October 5, 1960; made available for printing May 8, 1961.

R. OLDENBURGER and R. SRIDHAR are with Purdue University, Lafayette, Ind.

Research forming the basis of this paper was supported by a grant from the National Science Foundation of the United States Government.

..." as the latter value "goes to zero." The input was, at that time, assumed to be a biased sine wave. The theory developed by Oldenburger was unlike that published by Minorsky regarding the asynchronous quenching of systems described by nonlinear differential equations.<sup>3</sup>

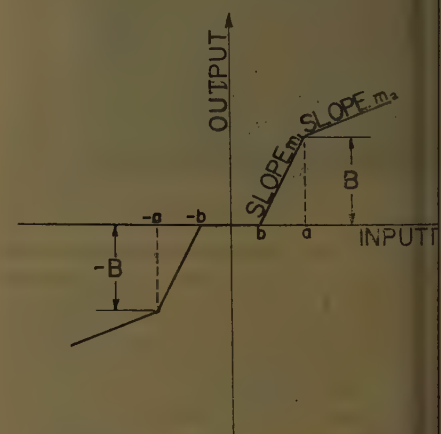


Fig. 1. Input-output characteristic, type nonlinearity



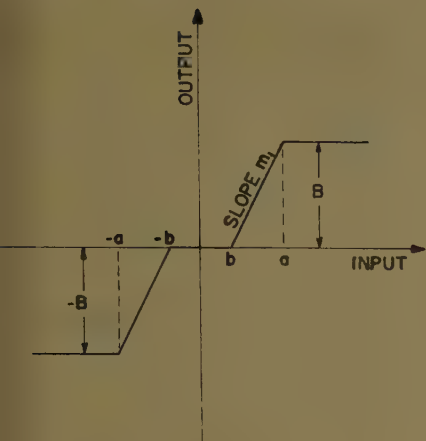


Fig. 2. Input-output characteristic, type B nonlinearity (limiter with dead band)

Oldenburger's signal-stabilization theory was extended, in collaboration with Nakada,<sup>4</sup> to self-oscillating systems; it is here further extended, in collaboration with Sridhar, to the case of a nonlinear system, which has one single-valued nonlinearity in the loop when the stabilizing signal belongs to a stationary random process with a Gaussian distribution and obeys the ergodic hypothesis. Now the definition of equivalent gain takes on wider scope than originally stated, and thereby becomes much more useful. The original actually turns out to be a particular case within the more general meaning of the new definition. Using this general precept, a method was found and is here presented for determining the mean-square value of noise to be injected at the input to certain self-oscillating nonlinear systems in order to reduce the output hunt to a desired value.

### Derivation of Equivalent Gains

Only nonmemory-type nonlinearities,<sup>5</sup> whose output  $y$  may be represented

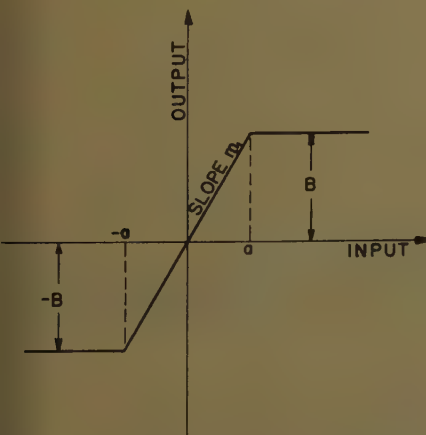


Fig. 3. Input-output characteristic, type C nonlinearity (limiter)

as a single-valued odd function  $f(x)$  of the input  $x$ , are within the restrictions set for this paper. The two general types to be discussed are those nonlinearities (1) whose input-output characteristic consists of straight-line segments, designated as "piecewise-linear," and (2) whose input-output characteristic may be described by a polynomial, designated as "polynomial."

The equivalent gain  $g(m)$  of a nonmemory-type of nonlinearity is the ratio between the average value  $A_v$  of the output  $y$  and the average value  $m$  of the input  $x$ , when the input to the nonlinearity consists of a stationary random process with a Gaussian distribution. The first probability density function  $p_1(x)$  of the input is given by

$$p_1(x) = \frac{1}{\sqrt{2\pi\phi_0}} \exp \left[ -\frac{(x-m)^2}{2\phi_0} \right] \quad (1)$$

where  $\phi_0$  is the variance of the input. Thus

$$g(m) = \frac{A_v}{m} \quad (2)$$

The limiting value of the equivalent gain defined by equation 2, when  $m \rightarrow 0$ , reduces to the equivalent gain defined by Oldenburger and Liu. This limiting value will be denoted as  $g(0)$ .

Equivalent gains of the two general types of nonlinearities will be derived next.

### Piecewise-Linear-Type Nonlinearity

The different piecewise-linear-type nonlinearities are shown in Figs. 1 through 7, and their input-output characteristics are described in Appendix I. Fig. 1 represents the most general type and the other six are special cases of this nonlinearity. They will be classified alphabetically from A to G for the sake of clarity. Thus, Fig. 1 corresponds to A, Fig. 2 to B, and so on. To distinguish between expressions for the equivalent gains of each nonlinearity, classifying letters of the alphabet will be added as subscripts to  $g(m)$ . Thus  $g_a(m)$  corresponds to the equivalent gain of the type-A nonlinearity, while  $g_a(0)$  corresponds to the gain of the same type when evaluated at  $m=0$ .

The only expression for equivalent gain which must be derived is that of the type-A nonlinearity. Expressions for other types are obtained by considering them as special cases of the type-A nonlinearity.

The general expression for  $A_v$  is

$$A_v = \int_{-\infty}^{\infty} f(x)p_1(x)dx \quad (3)$$

The expression for  $A_v$ , evaluated in

Appendix I for the type-A nonlinearity, yields

$$g_a(m) = \frac{m_1 - m_2}{m} \sqrt{\frac{\phi_0}{2\pi}} \left[ \exp \left\{ -\frac{(a+m)^2}{2\phi_0} \right\} - \exp \left\{ -\frac{(a-m)^2}{2\phi_0} \right\} \right] + m_2 - \frac{m_1}{m} \sqrt{\frac{\phi_0}{2\pi}} \left[ \exp \left\{ -\frac{(b+m)^2}{2\phi_0} \right\} - \exp \left\{ -\frac{(b-m)^2}{2\phi_0} \right\} \right] + \frac{m_1 - m_2}{2} \times \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) + \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) + \frac{m_1 - m_2}{m} \times \frac{a}{2} \times \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] - \frac{m_1}{2} \left[ \operatorname{erf} \left( \frac{b+m}{\sqrt{2\phi_0}} \right) + \operatorname{erf} \left( \frac{b-m}{\sqrt{2\phi_0}} \right) \right] - \frac{m_1}{m} \times \frac{b}{2} \times \left[ \operatorname{erf} \left( \frac{b+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{b-m}{\sqrt{2\phi_0}} \right) \right] \right] \quad (4)$$

Taking the limiting value of the right side of equation 4 and simplifying the resulting expression gives

$$g_a(0) = m_2 \left[ 1 - \operatorname{erf} \left( \frac{a}{\sqrt{2\phi_0}} \right) \right] + m_1 \left[ \operatorname{erf} \left( \frac{a}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{b}{\sqrt{2\phi_0}} \right) \right] \quad (5)$$

Expressions for equivalent gains of the other nonlinearities are now easily derived from equations 4 and 5, as illustrated in Appendix II.

Equivalent gains of the type-C nonlinearity (limiter) and the type-G nonlinearity (ideal relay) are shown in Figs. 8 and 9, respectively.

### Polynomial Type of Nonlinearity

The general polynomial type of nonlinearity is described by equation 6.

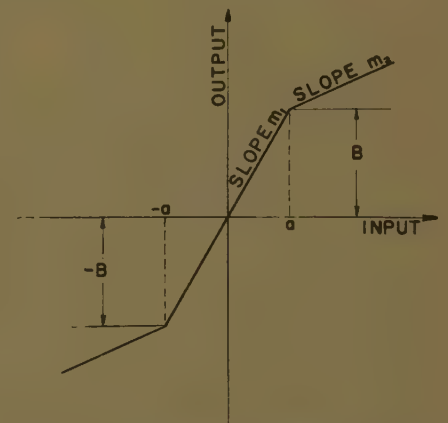


Fig. 4. Input-output characteristic, type D nonlinearity

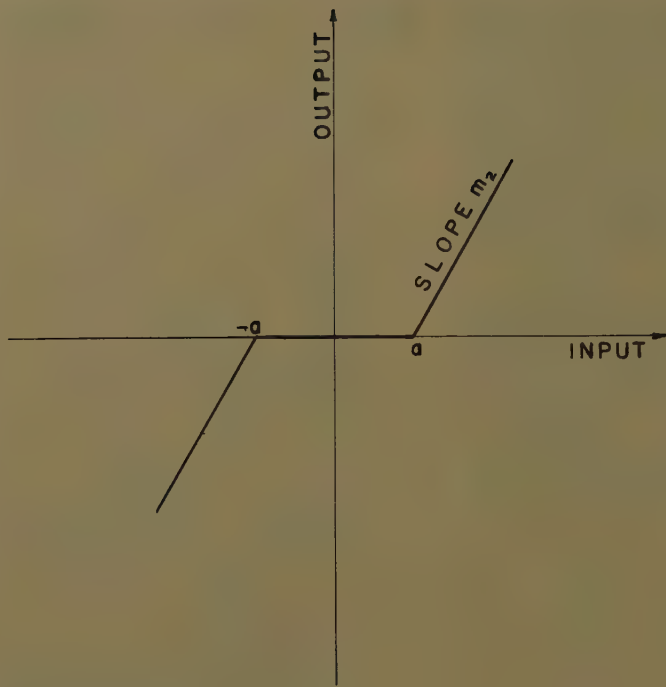


Fig. 5. Input-output characteristic, type E nonlinearity (dead-zone element)

There is no loss of generality in assuming that  $n$  is odd in

$$y = f(x) = c_n x^n + c_{n-1} x^{n-2} |x| + c_{n-2} x^{n-2} + \dots + c_{n-3} x^{n-4} |x| + \dots + c_2 x |x| + c_1 x \quad (6)$$

The absolute-value signs have been used in the right side of equation 6 to make  $f(x)$  an odd function of  $x$ .

In deriving the equivalent gain, each term of the polynomial is considered separately. Thus, for purpose of derivation, two types of nonlinearities are assumed. There is no loss in generality by this assumption since every term in the polynomial belongs to one of the two types shown in equations 7 and 8.

$$y = f(x) = cx^n \quad \text{for } n \text{ odd} \quad (7)$$

$$y = f(x) = \begin{cases} +cx^n & x \geq 0 \\ -cx^n & x < 0 \end{cases} \quad \text{for } n \text{ even} \quad (8)$$

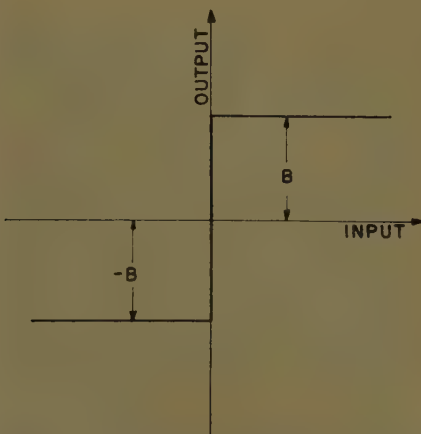


Fig. 7. Input-output characteristic, type G nonlinearity (ideal relay)

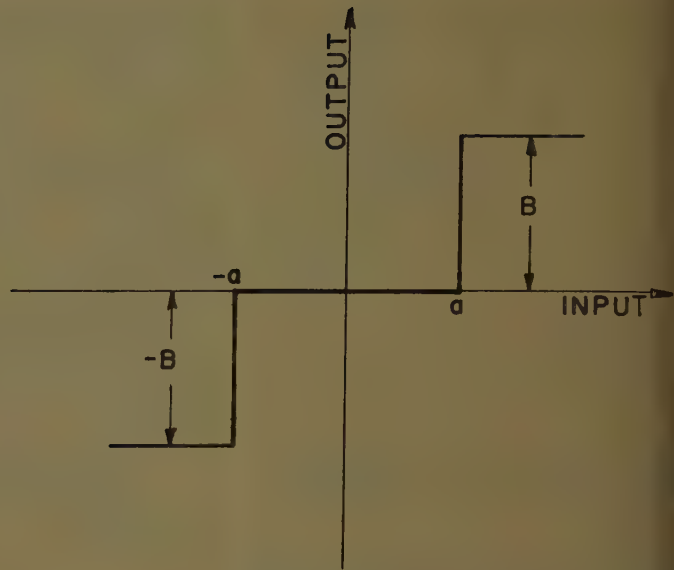


Fig. 6. Input-output characteristic, type F nonlinearity (relay with dead-band)

Also in Appendix III, the equivalent gain  $g_2(m)$  of the nonlinearity defined by equation 8 is given by

$$g_2(m) = \frac{A_v}{m} = \frac{c}{\sqrt{\pi}} m^{n-1} \left[ \sum_{k=0,2,4}^n \frac{n!}{(n-k)!} \times \frac{1}{k!} \times \left( \frac{\sqrt{2\phi_0}}{m} \right)^k \times 2 \int_0^{\frac{m}{\sqrt{2\phi_0}}} z^k \times \exp(-z^2) dz + \sum_{k=1,3,5}^{n-1} \frac{n!}{(n-k)!} \times \frac{1}{k!} \times \left( \frac{\sqrt{2\phi_0}}{m} \right)^k \times 2 \int_{\frac{m}{\sqrt{2\phi_0}}}^{\infty} z^k \exp(-z^2) dz \right] \quad (11)$$

The derivation in Appendix III shows that the nonlinearity corresponding to equation 7 has an equivalent gain  $g_1(m)$  of the form

$$g_1(m) = \frac{A_v}{m} = \frac{cm^{n-1}}{\sqrt{\pi}} \sum_{k=0,2,4}^{n-1} \times \frac{n!}{(n-k)!} \frac{1}{k!} \left( \frac{2\phi_0}{m} \right)^k \Gamma\left(\frac{k+1}{2}\right) \quad (9)$$

Taking the limiting value of the right side of equation 9 when  $m \rightarrow 0$  yields

$$g_1(0) = \frac{c}{\sqrt{\pi}} n (\sqrt{2\phi_0})^{n-1} \Gamma\left(\frac{n}{2}\right) \quad (10)$$

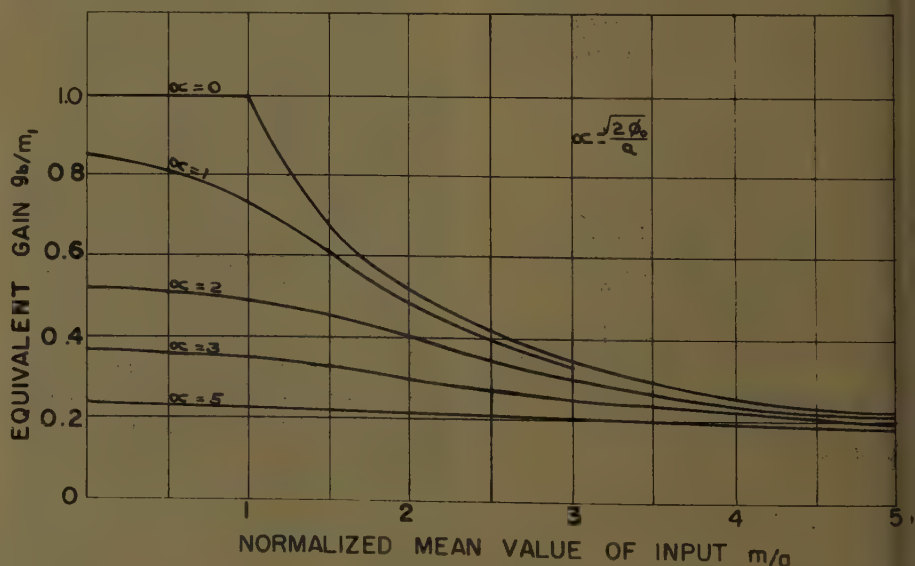


Fig. 8. Equivalent gain of limiter



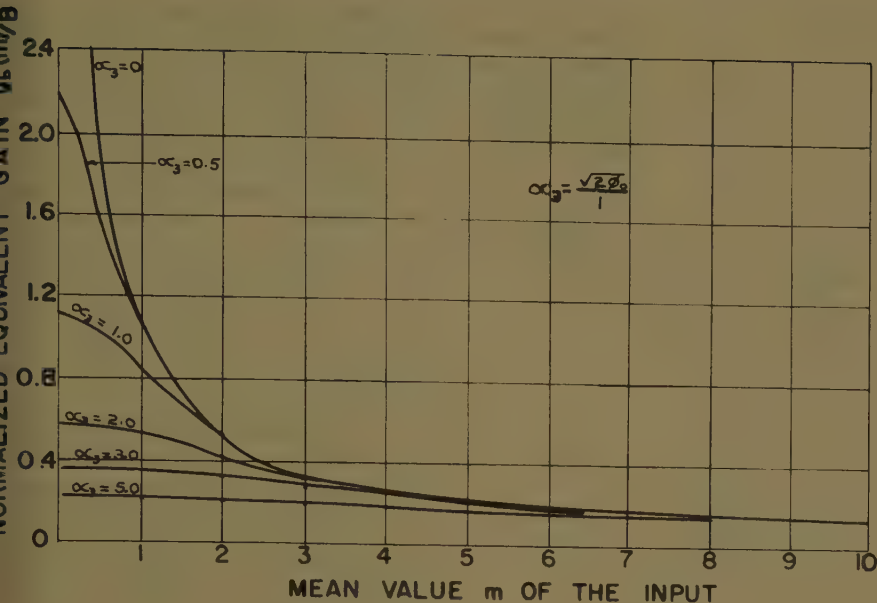


Fig. 9. Equivalent gain of ideal relay

Taking the limiting value of the right side of equation 11, when  $m \rightarrow 0$ , yields

$$g_2(0) = \frac{2c}{\sqrt{\pi}} \times n(\sqrt{2\phi_0})^{n-1} \int_0^\infty z^{n-1} \exp(-z^2) dz \quad (12)$$

The equivalent gain of the general polynomial-type nonlinearity, defined by equation 6, will equal the sum of all terms having the form seen in the right sides of equations 9 and 11.

The particular case when  $n$  equals zero in equation 8 results in the input-output characteristic of the ideal relay (type-G nonlinearity). Equating  $n$  to zero in the right side of equation 11 should yield the equivalent gain of the ideal relay which as derived in Appendix I. Thus

$$g_2(m) = \frac{c}{m} \operatorname{erf} \left( \frac{m}{\sqrt{2\phi_0}} \right) \text{ for } n=0 \quad (13)$$

The right sides of equations 13 and 30 are—as they should be—the same.

### Equivalent Admittance of a Nonlinearity

Before the theory of signal stabilization with random inputs is used in practical systems, it is advisable to understand “equivalent admittance” of a nonlinearity and also know how it may be determined.

Definition: The equivalent admittance  $J_{nl}(A, \phi_0)$  of a nonlinearity is the complex ratio of the amplitude  $P$  of the output component, which has the same frequency as the input

Table I. Fundamental and Third-Harmonic Components of Output Waveforms Shown in Figs. 10 and 11

$\sqrt{2\phi_0}$ ■ (Normalized Rms Value of Noise)	0	1	2	3
When $A/a$ (Normalized Amplitude of Sine Wave) equals 1				
$P_1/m_1a$ .....	1.0	0.745	0.507	0.3435
$P_3/m_1a$ .....	0	0.0323	0.01165	0.00445
When $A/a$ (Normalized Amplitude of Sine Wave) equals 2				
$P_1/m_1a$ .....	1.224	1.092	0.861	0.652
$P_3/m_1a$ .....	0	0.149	0.0538	0.0187

sine wave to the amplitude  $A$  of the input sine wave. (In this definition, it is assumed that the input consists of a sine wave, together with Gaussian noise, and that the power spectral density of the noise has no predominant component with the same frequency as the input sine wave.)

Notice the similarity between this definition and the describing function of a nonlinearity, which is  $J_{nl}(A, 0)$ .<sup>6</sup> For the types of nonlinearities considered in this paper, the equivalent admittance is a real quantity.

### Determination of Equivalent Admittance

Knowledge must be acquired regarding transmission of a sine wave through a nonlinearity in the presence of noise before equivalent admittance can be determined. The concept of equivalent gain may be used, as subsequently explained, to determine the amplitude of the output component of the nonlinearity having the same frequency as the input sine wave.

If the input to a nonlinearity consists

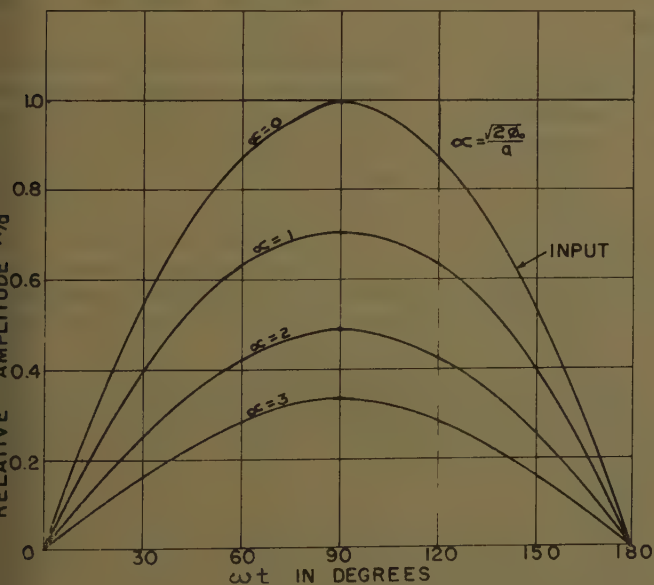


Fig. 10. Waveform, output of limiter for sinusoidal input  $A \sin \omega t$  (only half of period shown)

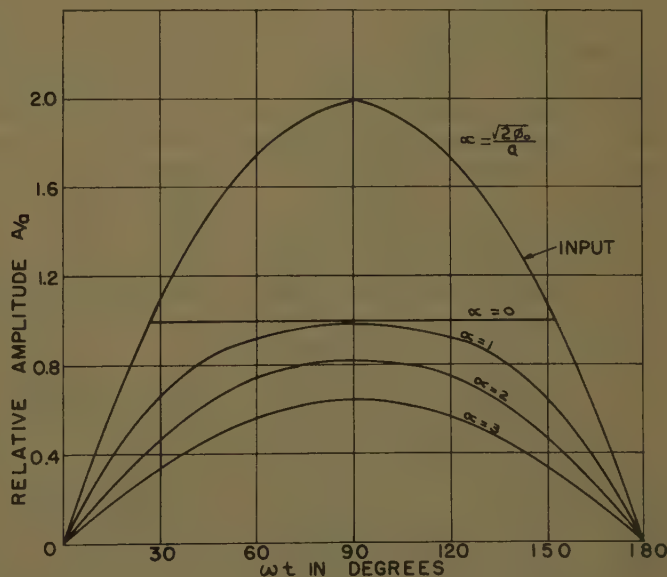


Fig. 11. Same as Fig. 10 for  $\frac{A}{a} = 2$  (only half of period shown)

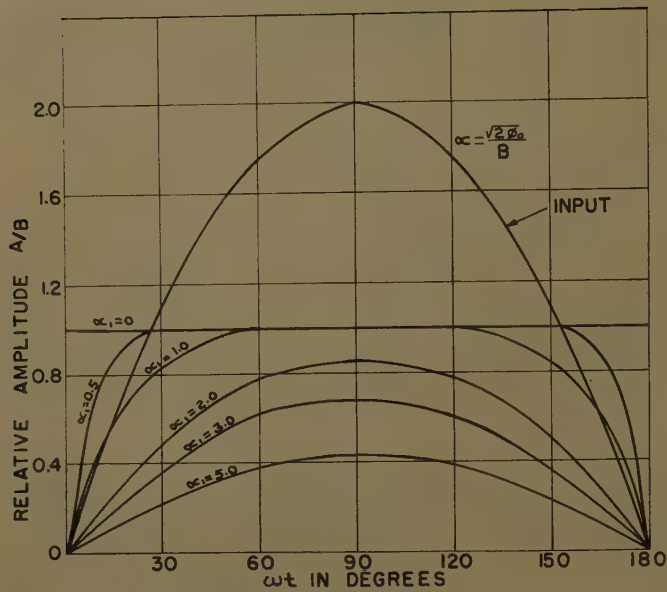


Fig. 12  
Waveform, output of  
ideal relay for sinu-  
soidal input  $A \sin$   
 $\omega t$  (only half of  
period shown)

of zero-mean noise, together with an additional signal of constant value, then the gain or amplification of the d-c component of the input is the equivalent gain. No restrictions are placed on the frequency components of the noise. Now assume that the additional signal is varying instead of constant. If the lowest frequency component of the spectral density of the noise is appreciably higher than the highest frequency component of the additional signal, a reasonable conclusion is that the additional signal is almost constant compared with the noise. The validity of this assumption increases with an increase in the lowest frequency component of the spectral density of the noise.

Hence, if the input to a nonlinearity consists of a stabilizing signal and an additional signal, the latter varying slowly in comparison with the former, then the instantaneous gain of the additional signal approximately equals the equivalent gain.

This concept will be used now to determine the gain of a sinusoidal signal through a nonlinearity when noise is present.

Let this noise be stationary and have Gaussian distribution with a mean of zero and a mean-square value  $\phi_0$  (note that the mean-square value is equal to the variance in this case), and let the sine wave be of amplitude  $A$  and frequency  $\omega$ . Consider half of one period of the sine wave as shown in Fig. 10. Divide one period into  $2n$  intervals, which is the same as dividing half a period into  $n$ -equal intervals, and erect ordinates at mid-points of each interval. When using the equivalent-gain technique, assume that the sine function is a constant, equal to the mid-ordinate value over each inter-

val. This is essentially equivalent to approximating the sine wave by a series of steps.

Consider the  $k$ th interval. The abscissa of its middle point is  $A \sin \pi/n(k - 1/2) = k_1$ . Since the input is assumed to be a constant equal to  $k_1$  over the  $k$ th interval, the output value of the nonlinearity over this interval is  $k_1 g(k_1)$ , where  $g(k_1)$  is the equivalent gain of the nonlinearity for an argument  $k_1$ . The value of the output over each  $2n$  interval may be determined similarly and the output waveform plotted. A series of steps comprise the output response. Fundamental and higher components of the output waveform may be determined now, using any of the well-known graphical or numerical methods.

Since the system beyond the nonlinearity is assumed to be essentially a low-pass filter, the components of the noise at the output may be neglected.

How to apply the equivalent-admittance method will be illustrated next by showing the transmission of a sine wave in the presence of noise through a limiter and a relay.

Figs. 10 and 11 show the input sine wave and the output waveforms for a limiter when  $A/a$  (normalized amplitude of the sine wave) is 1 and 2, respectively, and  $\sqrt{2\phi_0}/a$  (normalized rms value of

the noise) is 0, 1, 2, and 3. For convenience, a smooth curve has been drawn through the computed mid-ordinate points on the output waveform. One period of the sine wave has been divided into equal parts for this computation.

The fundamental and third-harmonic components of the output waveform in Figs. 10 and 11 were determined by using a 36-ordinate scheme over one full period, and the results are summarized in Table I, where  $P_1$  and  $P_3$  refer to the fundamental and third-harmonic components, respectively.

Similarly, Fig. 12 shows the input sine wave and the output waveforms for an ideal relay when  $A/B$  (normalized input amplitude) is 2 and  $\sqrt{2\phi_0}/B$  is 0.5, 1, 2, and 5. Fundamental and third-harmonic components of the output waveforms are shown in Table II. These components may be calculated in the same way for different nonlinearities by using knowledge of the equivalent gain of the nonlinearity. The ratio of this fundamental component amplitude to the input sine wave amplitude is the equivalent admittance. Equivalent-admittance curves for the limiter and the relay are shown in Figs. 13 and 14.

## Reduction of System Hunting to Desired Value

Signal stabilization is generally effective in a system having nonlinearity whose equivalent admittance  $J_{nl}(A, \phi_0)$  is monotonically nondecreasing in both  $A$  and  $\phi_0$ . The system should be absolutely stable below a particular value of gain when the nonlinearity is replaced by a simple gain, which means that the linear part of the system should not be conditionally stable.

In the discussion, it is assumed that the linear part of the system following the nonlinearity is such a good low-pass filter that noise feedback may be neglected.

In Fig. 15,  $G_1(s)$ ,  $G_2(s)$ , and  $H(s)$  represent the linear elements in the system and  $NL$  is the nonlinear element. The product  $G_1(s) \times G_2(s) \times H(s)$  is the linear part of the loop and is designated  $G(s)$ .

Table II. Fundamental and Third-Harmonic Components of Output Waveform Shown in Fig. 1

$\frac{\sqrt{2\phi_0}}{B}$ (Normalized Rms Value of Noise)	0	0.5	1	2	3	5
When $A/B$ (Normal Input Amplitude) equals 2						
$P_1/B$ .....	1.275	1.252	1.17	0.895	0.623	0.4
$P_3/B$ .....	0.3657	0.2087	0.0632	0.0207	0.0	0.0



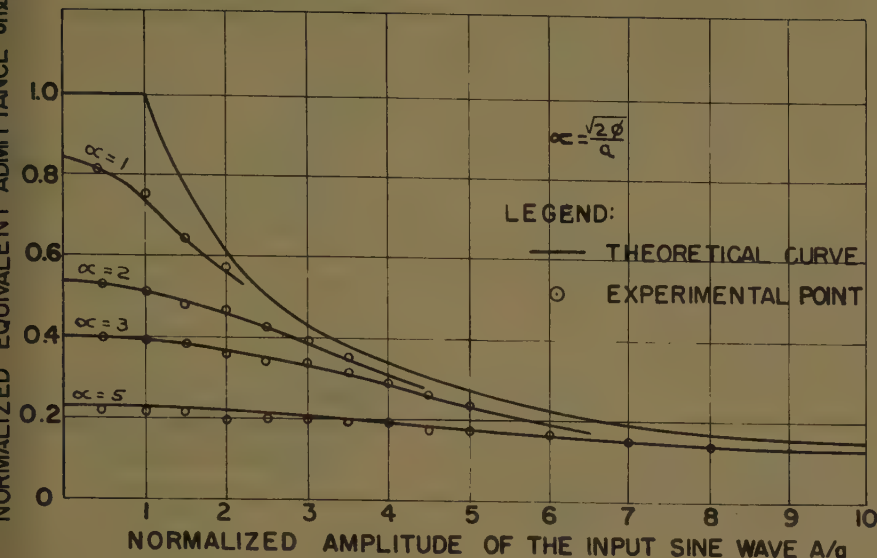


Fig. 13. Comparison of experimental and theoretical equivalent admittances of limiter

With no noise at the input to the nonlinearity, the amplitude and frequency of the system's output hunt, if any, may be determined by using the well-known describing-function technique. The intersections of the polar plot of  $G(j\omega)$  are investigated as is the negative of the describing-function's reciprocal on the complex plane of the nonlinearity. As should be remembered, the describing function is a special case of the equivalent admittance.

What will happen if Gaussian noise with a mean-square value of  $\phi_0$  is injected at the input to the nonlinearity? This noise is going to affect the characteristic of the nonlinearity and, hence, the amplitude of the output hunt. Using the same type of argument which justified replacing the nonlinearity with its describing function when the input was a pure sine wave, we can now justify replacing the nonlinearity with its equivalent admittance, since the input is a pure sine wave together with Gaussian noise.

In the latter case, the amplitude and frequency of hunt is determined by investigating the intersection of the polar plot of the linear part—system  $G(j\omega)$ —and the reciprocal's negative of the equivalent admittance. As the equivalent admittance of every nonlinearity considered in this paper is real and positive, the reciprocal's negative of the equivalent admittance will lie along the negative real axis in the complex plane. The frequency of auto-oscillations is therefore determined by the intersection of  $G(j\omega)$  with this axis. Since this point does not vary with changes in  $A$  or  $\phi_0$ , the frequency of auto-oscillations is not affected by the injection of noise. For a particular mean-square value of noise,

the amplitude of auto-oscillations may be determined from the intersection point.

From the standpoint of signal stabilization, however, our sole interest in the mean-square value of noise lies at the nonlinearity input, where hunting can be reduced to a desired value. Suppose the original hunt of the system is  $d_1$  units, and this is to be reduced to  $d_2$  units. A hunt amplitude of  $d_2$  units at the system output will result in an amplitude of  $A_2$  units at the input to the nonlinearity. This value is easily calculated when the actual transfer functions of the linear elements in the loop are known. The value of  $\phi_{01}$  now sought for is one which causes  $J_{n1}(A_1, \phi_0)$  to intersect  $G(j\omega)$ .

The procedure to follow is quite simple, as will be made clear after considering

an example, based on the nonlinear feedback system which is shown in Fig. 16. Using the describing function for the limiter, it can be shown that the amplitude and frequency of hunt of the system output are 63.5 units and 10 radians per second, respectively. As equivalent-admittance curves for the limiter are shown in Fig. 13, the output hunt is seen to be 58 units and 32 units for  $\sqrt{2\phi_0}/a = 3$  and 5, respectively. Assuming that the noise fed back is negligible, then the mean-square value of noise at the limiter input is the same as at the system input.

## Verification of Equivalent Admittance

An analog computer was used for experimentally verifying the equivalent admittance of the limiter. This verification should also serve to verify the equivalent gain, since this gain was used extensively in deriving equivalent-admittance curves.

The output of a commercial white-noise generator was passed through a second-order shaping filter and was added to the output of a sine-wave generator. The output of the adder was thus the sum of a sine wave and essentially Gaussian noise, which sum was the input to a limiter circuit. The fundamental component of the limiter output was measured with a sharply tuned filter. From this, the equivalent admittance of the limiter was evaluated. The experimentally determined points appear as dots in Fig. 13, which shows close agreement between theoretically determined curves and these points.

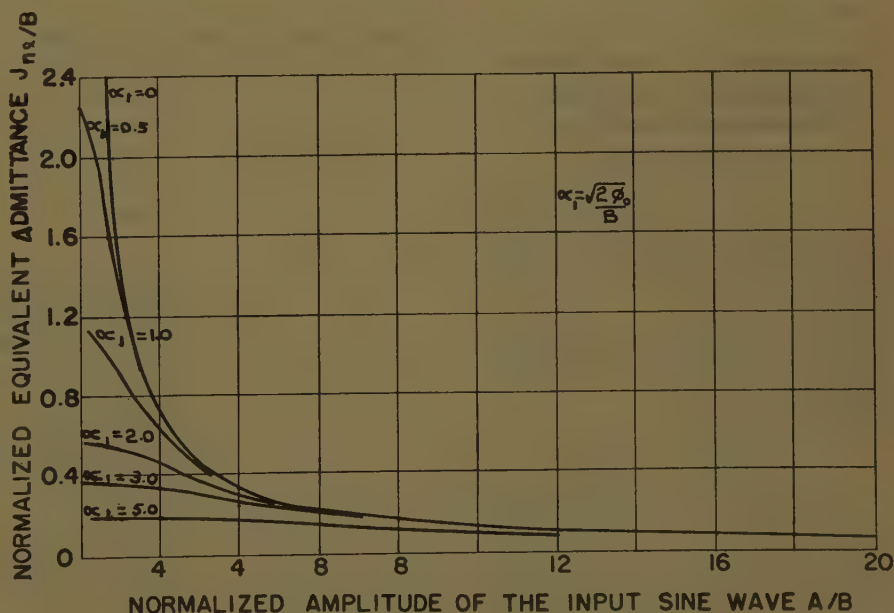


Fig. 14. Equivalent admittance  $J_{n1}(A, \phi_0)$  of ideal relay

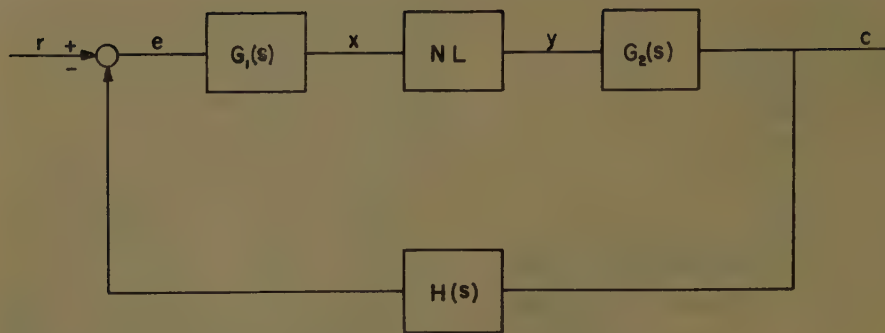


Fig. 15. General type of nonlinear feedback system

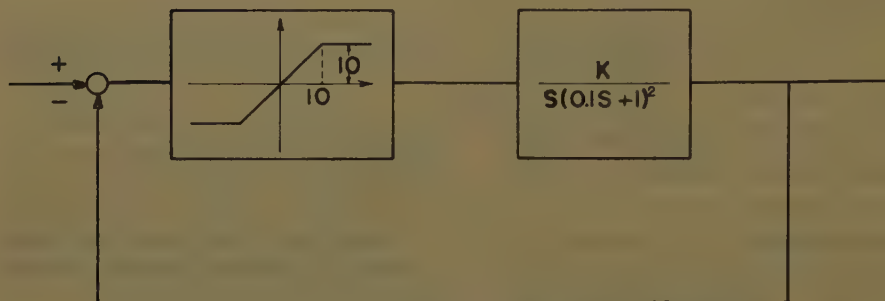


Fig. 16. A particular nonlinear feedback system

### Amplitude of Auto-Oscillations of a Closed-Loop System

The closed-loop system in Fig. 16 was simulated on an analog computer. Noise with different mean-square values was injected at the system input for various values of the gain constant  $K$ , and amplitudes of the output hunt were measured. Results of this experiment, together with the theoretically calculated points, are shown in Fig. 17. Close agreement is again evident between the theory and experiment.

### Appendix I. Input-Output Characteristics of Certain Piecewise-Linear-Type Nonlinearities

Input-output characteristics of the piecewise-linear-type nonlinearities shown in Figs. 1 through 7 are here described as a function. The alphabetical classification of each nonlinearity is also included.

1. The input-output characteristic of type A nonlinearity, illustrated in Fig. 1, is described by equation 14

$$y=f(x)=\begin{cases} m_2(x-a)+B & x \geq a \\ m_1(x-b) & b < x < a \\ 0 & -b \leq x \leq b \\ m_1(x+b) & -a < x < -b \\ m_2(x+a)-B & x \leq -a \end{cases} \quad (14)$$

2. The input-output characteristic of type B nonlinearity (limiter with dead band), seen in Fig. 2, is described by the following equation:

$$y=f(x)=\begin{cases} B & x \geq a \\ m_1(x-b) & b < x < a \\ 0 & -b \leq x \leq b \\ m_1(x+b) & -a < x < -b \\ -B & x \leq -a \end{cases} \quad (15)$$

3. The input-output characteristic of type C nonlinearity (limiter), Fig. 3, is

$$y=f(x)=\begin{cases} B & x \geq a \\ m_1x & |x| < a \\ -B & x \leq -a \end{cases} \quad (16)$$

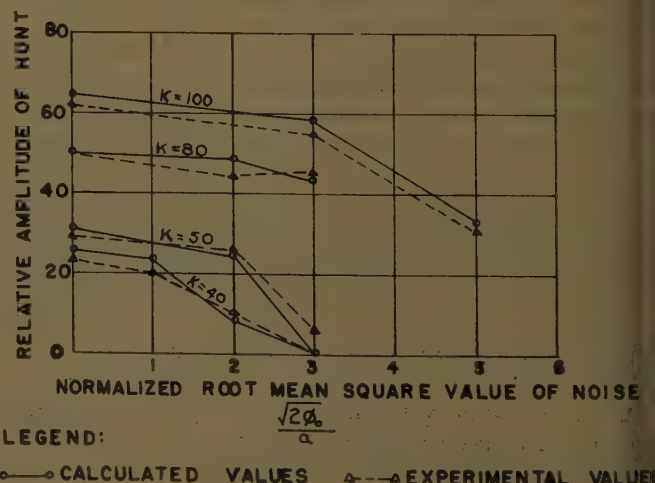
4. Type D nonlinearity, shown in Fig. 4, has the input-output characteristic

$$y=f(x)=\begin{cases} m_2(x-a)+B & x \geq a \\ m_1x & |x| < a \\ m_2(x+a)-B & x \leq -a \end{cases} \quad (17)$$

5. Fig. 5 shows the input-output characteristic of type E nonlinearity described by

$$y=f(x)=\begin{cases} m_2(x-a) & x \geq a \\ 0 & |x| < a \\ m_2(x+a) & x \leq -a \end{cases} \quad (18)$$

Fig. 17. Amplitude, oscillation of system output, Fig. 16



6. The input-output characteristic of type F nonlinearity (relay with dead band), Fig. 6, is

$$y=f(x)=\begin{cases} +B & x \geq a \\ 0 & |x| < a \\ -B & x \leq -a \end{cases} \quad (19)$$

7. The ideal relay or type G nonlinearity has the input-output characteristic shown in Fig. 7 and expressed by

$$y=f(x)=\begin{cases} +B & x > 0 \\ 0 & x = 0 \\ -B & x < 0 \end{cases} \quad (20)$$

### Appendix II. Expression for Equivalent Gains of Some Piecewise-Linear-Type Nonlinearities

Expressions for equivalent gains of type B through G nonlinearities are derived by considering them as modifications of type A. By properly modifying the right sides of equations 4 and 5, expressions are found for equivalent gains of required nonlinearities, each of which is considered separately.

1. Type B nonlinearity (limiter with dead band) is reduced from type A when  $m_2=0$ . A. If  $m_2=0$  in the right sides of equations 4 and 5, the result is

$$g_b(m) = \frac{m_1}{m} \sqrt{\frac{\phi_0}{2\pi}} \left[ \exp \left\{ -\frac{(a+m)^2}{2\phi_0} \right\} - \exp \left\{ -\frac{(a-m)^2}{2\phi_0} \right\} - \exp \left\{ -\frac{(b+m)^2}{2\phi_0} \right\} + \exp \left\{ -\frac{(b-m)^2}{2\phi_0} \right\} \right] + \frac{m_1}{2} \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) + \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] + \frac{m_1}{2m} \left[ a \left\{ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right\} - b \left\{ \operatorname{erf} \left( \frac{b+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{b-m}{\sqrt{2\phi_0}} \right) \right\} \right] - \frac{m_1}{2} \left[ \operatorname{erf} \left( \frac{b+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{b-m}{\sqrt{2\phi_0}} \right) \right] \quad (21)$$



$$=m_1 \left[ \operatorname{erf} \left( \frac{a}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{b}{\sqrt{2\phi_0}} \right) \right] \quad (22)$$

Type B nonlinearity reduces to the type limiter when  $b=0$  in the former. Let  $b \rightarrow 0$  in the right sides of equations 21 and 22 yields

$$\begin{aligned} &= \frac{m_1}{m} \times \sqrt{\frac{\phi_0}{2\pi}} \left[ \exp \left\{ -\frac{(a+m)^2}{2\phi_0} \right\} - \exp \left\{ -\frac{(a-m)^2}{2\phi_0} \right\} \right] + \frac{m_1}{2} \times \\ &\left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) + \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] + \\ &\frac{m_1}{m} \times \frac{a}{2} \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] \quad (23) \end{aligned}$$

$$=m_1 \operatorname{erf} \left( \frac{a}{\sqrt{2\phi_0}} \right) \quad (24)$$

Type A nonlinearity reduces to type D when  $b=0$  in the former, yielding equation 25 when  $b \rightarrow 0$  in equations 4 and 5

$$\begin{aligned} &= \frac{m_1 - m_2}{m} \sqrt{\frac{\phi_0}{2\pi}} \left[ \exp \left\{ -\frac{(a+m)^2}{2\phi_0} \right\} - \exp \left\{ -\frac{(a-m)^2}{2\phi_0} \right\} \right] + m_2 + \frac{m_1 - m_2}{2} \times \\ &\left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) + \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] + \\ &\frac{m_1 - m_2}{m} \times \frac{a}{2} \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] \quad (25) \end{aligned}$$

equation 26

$$\begin{aligned} &= m_2 \left[ 1 - \operatorname{erf} \left( \frac{a}{\sqrt{2\phi_0}} \right) \right] + \\ &m_1 \left[ \operatorname{erf} \left( \frac{a}{\sqrt{2\phi_0}} \right) \right] \quad (26) \end{aligned}$$

Type D nonlinearity reduces to the type limiter when  $m_1=0$  in the former and when  $m_1 \rightarrow 0$  in the right sides of equations 25 and 26

$$\begin{aligned} &= -\frac{m_2}{m} \sqrt{\frac{\phi_0}{2\pi}} \left[ \exp \left\{ -\frac{(a+m)^2}{2\phi_0} \right\} - \exp \left\{ -\frac{(a-m)^2}{2\phi_0} \right\} \right] + m_2 - \\ &\frac{m_2}{m} \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) + \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] - \\ &\frac{m_2}{m} \times \frac{a}{2} \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] \quad (27) \end{aligned}$$

$$=m_2 \left[ 1 - \operatorname{erf} \left( \frac{a}{\sqrt{2\phi_0}} \right) \right]$$

5. The expression for equivalent gain for type F nonlinearity (relay with dead band) can be derived more easily from the original definition than by considering it as a modification of any other type. Substituting the expression for  $f(x)$  from equation 19 for  $A$ , in equation 3 and then substituting the resulting expression for  $A$ , in equation 2 yields the following

$$g_f(m) = \frac{B}{2m} \left[ \operatorname{erf} \left( \frac{a+m}{\sqrt{2\phi_0}} \right) - \operatorname{erf} \left( \frac{a-m}{\sqrt{2\phi_0}} \right) \right] \quad (28)$$

and

$$g_f(0) = B \sqrt{\frac{2}{\phi_0 \pi}} \times \exp \left( -\frac{a^2}{2\phi_0} \right) \quad (29)$$

6. Now type F (ideal relay) reduces to the type G nonlinearity when  $a=0$  in the former. Letting  $a=0$  in the right side of equations 28 and 29 yields

$$g_g(m) = \frac{B}{m} \operatorname{erf} \left( \frac{m}{\sqrt{2\phi_0}} \right) \quad (30)$$

and

$$g_g(0) = B \sqrt{\frac{2}{\pi \phi_0}} \quad (31)$$

### Appendix III. Average Output Value of Polynomial-Type Nonlinearity

In this derivation, rather than consider the general polynomial-type nonlinearity described by equation 6, the two special cases described by equations 7 and 8 are viewed separately.

For case 1,  $n$  odd, substitutions from equations 1 and 7 in equation 3 will produce the following

$$A_g = \frac{c}{\sqrt{2\pi\phi_0}} \int_{-\infty}^{\infty} x^n \exp \left\{ -\frac{(x-m)^2}{2\phi_0} \right\} dx \quad (32)$$

By suitable change of variables and manipulation of the results, it can be shown that equation 32 yields

$$A_g = \frac{cm^n}{\sqrt{\pi}} \sum_{k=0,2,4}^{n-1} \frac{n!}{(n-k)!} \frac{1}{k!} \left( \frac{\sqrt{2\phi_0}}{m} \right)^k \gamma \times \left( \frac{k+1}{2} \right) \quad (33)$$

where  $\Gamma$  represents the gamma function.

For case 2,  $n$  even, substitutions from equations 1 and 8 in equations 3 will produce

$$\begin{aligned} A_g &= \frac{-c}{\sqrt{2\pi\phi_0}} \int_{-\infty}^0 x^n \exp \left\{ -\frac{(x-m)^2}{2\phi_0} \right\} dx + \\ &\frac{c}{\sqrt{2\pi\phi_0}} \int_0^{\infty} x^n \exp \left\{ -\frac{(x-m)^2}{2\phi_0} \right\} dx \quad (34) \end{aligned}$$

By suitable changes of variables and manipulation of the results, equation 34 yields

$$\begin{aligned} A_g &= \frac{-c}{\sqrt{\pi}} m^n \left[ \sum_{k=0,2,4}^n \frac{n!}{(n-k)!} \frac{1}{k!} \left( \frac{\sqrt{2\phi_0}}{m} \right)^k \times \right. \\ &2 \int_0^{\sqrt{2\phi_0}} z^k \exp(-z^2) dz + \\ &\sum_{k=1,3,5}^n \frac{n!}{(n-k)!} \frac{1}{k!} \left( \frac{\sqrt{2\phi_0}}{m} \right)^k \times \\ &\left. 2 \int_{\frac{m}{\sqrt{2\phi_0}}}^{\infty} z \exp(-z^2) dz \right] \quad (35) \end{aligned}$$

### References

1. SIGNAL STABILIZATION OF A CONTROL SYSTEM, R. Oldenburger. *Transactions, American Society of Mechanical Engineers*, New York, N. Y., vol. 79, no. 8, Nov. 1957, pp. 1869-72.
2. SIGNAL STABILIZATION OF A CONTROL SYSTEM, R. Oldenburger, C. C. Liu. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 78, May 1959, pp. 96-100.
3. ON ASYNCHRONOUS ACTION, N. Minorsky. *Journal, Franklin Institute, Philadelphia, Pa.*, vol. 259, no. 3, Mar. 1955, pp. 209-19.
4. SIGNAL STABILIZATION OF SELF-OSCILLATING SYSTEMS, R. Oldenburger, T. Nakada. *Transactions, Professional Group on Automatic Control, Institute of Radio Engineers*, New York, N. Y., vol. AC-6, no. 3, Sept. 1961.
5. A GENERAL METHOD OF DERIVING THE DESCRIBING FUNCTION OF A CERTAIN CLASS OF NON-LINEARITIES, R. Sridhar. *Ibid.*, vol. AC-5, no. 2, June 1960, pp. 135-41.
6. ON SOME NONLINEAR PHENOMENA IN REGULATORY SYSTEMS (in Russian), L. C. Goldfarb. *Automatika i Telemekhanika*, Moscow, USSR, vol. 8, Sept.-Oct. 1947, pp. 349-83; also in *FREQUENCY RESPONSE* (book, in English), R. Oldenburger, editor. The Macmillan Company, New York, N. Y., 1956, pp. 239-59.

### Discussion

R. L. Cosgriff (Ohio State University, Columbus, Ohio): The authors are to be congratulated for focusing attention on characteristics of nonlinear systems excited by random signals. This discussion will review similar work and deal with extensions of the method presented in this paper.

A classic in the control field, which tells how to linearize the characteristics of a nonlinear device by a secondary or dither signal, is reference 1, and a general linearizing process is described by Loeb in reference 2.

The problem of random signals discussed in this paper is so closely associated with the gain function defined by Booton<sup>3</sup> that their small signal gain  $g_1(0)$  is identical to Booton's gain function. Much of the paper's development parallels that of reference 4, although the techniques differ for obtaining the deterministic portion of the output signal.

Both treatments produce identical results for nonlinearities of the form  $y=x^n$  when  $n$  is odd. Reference 4 also gives the series form for the terms when  $n$  is even, and allows determination of the output's expected value for nonlinear blocks, characterized by more general types of nonlinearity

than were possible with the author's development. For example, the expected value of  $y$  for

$$y = x(dx/dt)^2$$

and for

$$x = m(t) + n(t)$$

where  $m(t)$  is a known signal and  $n(t)$  is random with mean zero is determined by first substituting for  $x$ , giving

$$y = [m(t) + n(t)] [(dm/dt)^2 + 2 \frac{dm}{dt} \frac{dn}{dt} + (dn/dt)^2]$$

The expected values of  $y$  become

$$y_e = m(t) \left[ \sigma_z^2 + \left( \frac{dm(t)}{dt} \right)^2 \right]$$

Here  $\sigma_z^2$  is the variance of  $dn/dt$ . For  $m(t) = A \cos \omega t$ , the value of  $y$ 's fundamental component is

$$\left( \sigma_z^2 + \frac{\omega^2}{4} \right) \cos \omega t$$

The authors' development of the gain

function for certain nonlinear elements, for sinusoidal input signals, and for a random signal—based upon the concepts of Kochenberger and others—should prove to be useful; see Figs. 13 and 14. I regret that curves for relay with dead band were not included. However, the omission is understandable when considering the difficulty of presenting the data alone.

#### REFERENCES

1. FUNDAMENTAL THEORY OF SERVOMECHANISMS (book), L. A. McColl. D. Van Nostrand Company, Inc., Princeton, N. J., 1945, pp. 79-87.
2. A GENERAL LINEARIZING PROCESS FOR NONLINEAR CONTROL SYSTEMS, J. M. Loeb. AUTOMATIC AND MANUAL CONTROL (book). Academic Press, Inc., New York, N. Y., 1952, pp. 275-83.
3. ANALYSIS OF NONLINEAR CONTROL SYSTEMS WITH RANDOM INPUTS, R. C. Booton. *Proceedings, Nonlinear Circuit Analysis Symposium*, April 24, 1953, pp. 369-92.
4. NONLINEAR CONTROL SYSTEMS (book), R. L. Cosgriff. McGraw-Hill Book Company, Inc., New York, N. Y., 1958, pp. 270-80.

R. Oldenburger and R. Sridhar: Prof. Cosgriff's discussion adds greatly to the value of our paper. Certain developments cited by him parallel part of our work, although

our main interest centers around stabilizing a feedback control system in limit cycle operation by injecting a random signal. This specific problem is not satisfactorily explained in any of the references mentioned in the discussion.

The theory of signal stabilization can be extended rather easily to nonlinearities more general than the ones considered. This is evidenced by Professor Cosgriff's example of the nonlinearity whose output involves the input and its derivative. However, this is a matter of details rather than generalities. The theory's extension to systems with nonsingle valued and other complicated nonlinearities warrants further work. A simple example is a system which has an element with hysteresis.

The theory developed to explain signal stabilization with random inputs evidently may be extended to yield more information. It can, for example, obtain stability information for nonlinear systems with a noise pickup in the loop. (See reference 1).

#### REFERENCE

1. STABILITY OF A NONLINEAR FEEDBACK SYSTEM IN THE PRESENCE OF GAUSSIAN NOISE, R. Sridhar, R. Oldenburger. *Paper no. 61-JAC-5*, American Society of Mechanical Engineers, New York, N. Y., 1961.

# A Graphical Method for Finding the Frequency Response of Nonlinear Closed-Loop Systems

A. S. McALLISTER  
STUDENT MEMBER AIEE

SEVERAL methods<sup>1-7</sup> have been developed for finding the steady-state sinusoidal response of closed-loop systems containing nonlinear elements. Most of these methods require either unnecessarily severe assumptions or long iterative procedures. This paper presents a graphical procedure for finding the frequency response of a large class of systems for which the describing function provides an adequate representation for the nonlinear elements. Attention is restricted to frequency-insensitive zero-memory nonlinear elements. The method is direct, requiring only an understanding of the point of view, and it can be extended to cover systems containing more than one nonlinear element.

## A Determination of Closed-Loop Sinusoidal Response

Consider the single-loop nonlinear system shown in the block diagram of Fig. 1. The blocks marked  $g_1$ ,  $g_2$ , and  $h$  are linear

elements whose input-output relations can be expressed by the usual transfer functions using the Laplace variables,  $G_1(s)$ ,  $G_2(s)$ , and  $H(s)$ . The block marked  $n$  is a nonlinear element whose instantaneous output can be expressed as a single-valued odd function of its instantaneous input, that is, this block is a frequency-independent zero-memory nonlinear element. Under these conditions, if the system input  $r(t)$  is a steady sinusoid, the other system signals  $e$ ,  $x$ ,  $y$ ,  $c$ , and  $f$  are generally periodic but not necessarily sinusoidal. In an important class of control systems the linear elements have low-pass characteristics, so that even though  $y$  may be non-sinusoidal, its higher harmonics are filtered as the signal passes around the loop, making  $x$  very nearly sinusoidal. Under these conditions only the fundamental component of  $y$  is of importance in analyzing the closed-loop sinusoidal behavior. Thus, with the limitations of frequency insensitivity and zero memory, the nonlinear element can be approximated as a

linear gain element whose gain  $N(X)$  is a function of only the amplitude of the input  $X$ . This gain is the describing function of the nonlinear element. Throughout this analysis it is assumed that the input of the nonlinear element is sinusoidal and the describing function is an adequate representation of the nonlinearity. An analysis of the system can be carried out if all the system signals are considered to be sinusoids of the same frequency with the input  $r$  is sinusoidal.

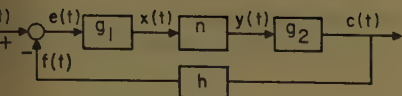
It should be noted that even if the system loop has a sufficiently low-pass characteristic, so that the describing function method may be used for sinusoidal analysis, the system output  $c$  need not necessarily be a pure sinusoid but can contain harmonics. If the values of any harmonic components of the output are desired, they can be approximately determined by first assuming that the describing function gives the correct value of  $N$ . This value is used to find the harmonics in  $y$ . The harmonics in  $c$  can then be found by applying  $G_2(j\omega)$  with the appropriate values of  $\omega$ .

Now, if all the signals of the system are sinusoidal, then the system is linear. For any given value of the gain,  $N$ , and

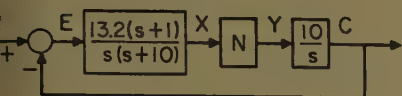
Paper 61-710, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department. Presented at the AIEE-AICHE-ASME-IREE Joint Automatic Control Conference, Boulder, Colo., June 28-30, 1961. Manuscript submitted September 8, 1960; made available for printing May 9, 1961.

A. S. McALLISTER is an Associate Professor, San Jose State College, San Jose, Calif.

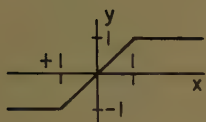




1. Feedback control system containing a nonlinear element



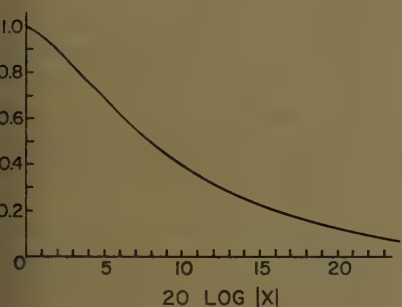
(A)



(B)

2. Specific example of nonlinear system

A—Block diagram  
B—Idealized saturation



3. Describing function for idealized saturation of Fig. 2

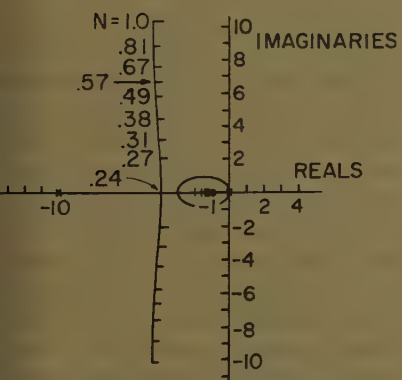


Fig. 4. Root locus for system of Fig. 2

Amplitude ratio and relative phase of  $C/R$  can be obtained in the usual manner.

$$\frac{C}{R} = \frac{G_1(j\omega)N(X)G_2(j\omega)}{1 + G_1(j\omega)N(X)G_2(j\omega)H(j\omega)} \quad (1)$$

Thus, the ratio  $C/R$  can be determined for a given  $N$  and frequency  $\omega$ . Un-

Fig. 5. Closed-loop frequency response of system of Fig. 2 for various values of  $N$

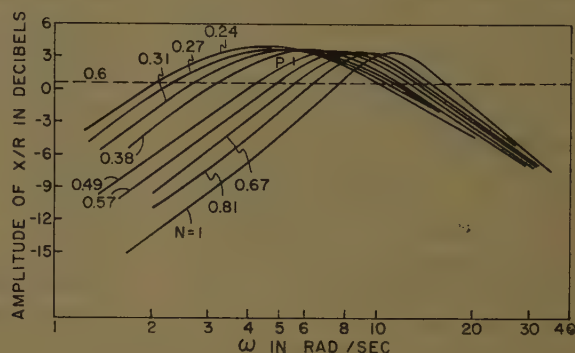
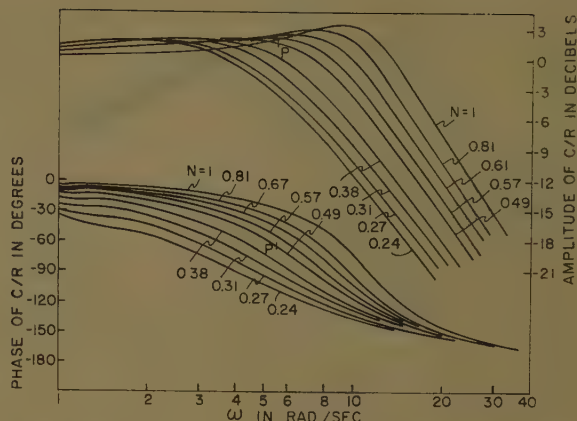


Fig. 6.  $X/R$  curves for system of Fig. 2

fortunately this is not the information that is desired; what is needed is the value of this ratio as a function of the amplitude and frequency of the input  $R$ . But the value of  $X$  (and therefore the value of  $N$ ) is not immediately obtainable when  $R$  is given. Indeed,  $X$  is a very complicated function of  $R$ . The transfer function relating the two,

$$\frac{X}{R} = \frac{G_1(j\omega)}{1 + G_1(j\omega)N(X)G_2(j\omega)H(j\omega)} \quad (2)$$

is a function of  $X$ . The solution of equation 2 for  $X$ , in terms of  $R$ , is usually prohibitively difficult even for the simplest nonlinearities.

However, if the point of view is slightly changed, a method which readily lends itself to graphical solution is presented. The solution of equation 2 is extremely difficult or sometimes impossible, when  $R$  is treated as the independent variable and the corresponding values of  $X$  and  $N$  are sought; however, the solution is simplified if  $X$  is taken as the independent variable. For a given amplitude of  $X$ ,  $N$  is a constant and  $X/R$  is an ordinary linear closed-loop transfer function which can be solved as a function of frequency  $\omega$  in any of the usual ways.<sup>1,8,9</sup> Thus,  $R$  can be found as a function of frequency for the given  $X$ . If this is done for a sufficient number of values of  $X$ , the difficulty will have been circumvented.  $R$  will be known for every frequency,  $\omega$ , and every gain value  $N$ . Now, when the amplitude and frequency of an input,  $R$ , are given,

the accumulated data are examined and the corresponding value of  $N$  is found. This value in equation 1 gives the desired output  $C$ .

It is desirable to have a complete set of data relating  $C$  to  $R$  for all amplitudes and frequencies of  $R$ ; this may be accomplished by utilizing all the data that relate  $N$  to  $R$  in the repeated solutions of equation 1. At first glance this appears to be a tremendous task, but with systematic procedures and graphical techniques it is no more difficult than the plotting of the closed-loop frequency response of a linear system for several values of loop gain. The following section illustrates such a graphical method of solution.

### Example Using a Graphical Approach

A specific example of the system of Fig. 1 is shown in Fig. 2(A) where  $n$  is taken as idealized saturation as shown in Fig. 2(B). The describing function,  $N$ , for saturation has been given by many authors<sup>10</sup> and is shown in Fig. 3. Though much of the work of this section is carried out in terms of the complex Laplace variables, it must be remembered that under the assumption of the describing function the analysis is only valid when  $x$  is less than the saturation value 1 or when  $s$  is purely imaginary,  $s = j\omega$ , that is, when the signals are all steady-state sinusoids of the same frequency.

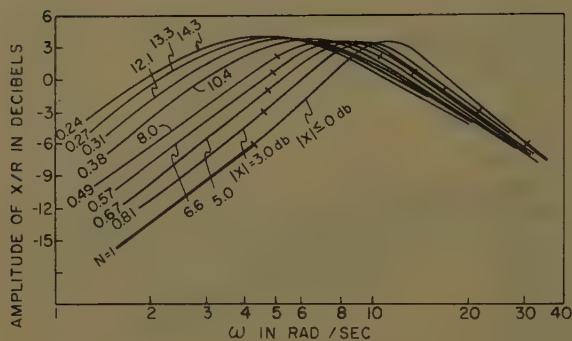


Fig. 7. Replot of Fig. 6 with  $|R| = 6$  db points identified

The two transfer functions which are of interest can now be written

$$\frac{C}{R} = \frac{13.2(s+1) N \frac{10}{s}}{1 + \frac{13.2(s+1) N \frac{10}{s}}{s^2(s+10)}} = \frac{N132(s+1)}{s^2(s+10) + N132(s-1)} \quad (3)$$

$$\frac{X}{R} = \frac{13.2(s+1)}{s(s+10)} \frac{1}{1 + \frac{13.2(s+1) N \frac{10}{s}}{s^2(s+10)}} = \frac{132s(s+1)}{s^2(s+10) + N132(s+1)} = \frac{C}{R} \frac{s}{N10} \quad (4)$$

To plot curves of these functions for  $s=j\omega$ , the roots of the denominator common to both functions must be determined. The location of the roots can be obtained as a function of  $N$  by means of the root locus<sup>1,11</sup>. The root locus for the system under consideration is shown in Fig. 4. The linear gain 132, is selected so that the complex roots have a damping factor,  $\zeta = 0.4$ , when operating in the linear region. The values of the describing function,  $N$ , are marked on the complex roots in Fig. 4. The real roots are not marked. The data from Fig. 4 are given in Table I, where

$$s^2(s+10) + 132N(s+1) = (s+p_3)(s+2\zeta\omega_n s + \omega_n^2) \quad (5)$$

132 $N$  is the loop gain,  $N$  the describing function,  $\zeta$  the damping factor of the

complex pair of roots,  $\omega_n$  the undamped natural frequency of the complex roots,  $-p_3$  the location of the third root, and  $|X|$  the value of input to the nonlinearity corresponding to  $N$  as obtained from Fig. 3. It can be seen in Table I that the values of gain were chosen so that the values of the damping factor,  $\zeta$ , would be simple. This was done so that standard templates could be used to draw the frequency response curves.

For constant values of  $N$ , all the poles and zeros of equations 3 and 4 are known, and frequency response curves can be drawn by plotting log-amplitude and phase versus log-frequency. Asymptotic approximations to these curves are not sufficient because these closed-loop functions involve complex poles and because fairly accurate frequencies at which certain amplitudes occur are wanted. However, if templates or similar devices are used, it is reasonably easy to add the effects of the various poles and zeros graphically to obtain the final curves. The curves used in this paper were drawn with the aid of a set of normalized curves which were placed under the final curve sheet and added as they were traced. The curves of equations 3 and 4 are shown in Figs. 5 and 6. Since only the amplitude of  $X$  is needed for determining the value of the describing function, the phase of  $X/R$  is not necessary and therefore has not been drawn.

In drawing Fig. 6 a simplification was

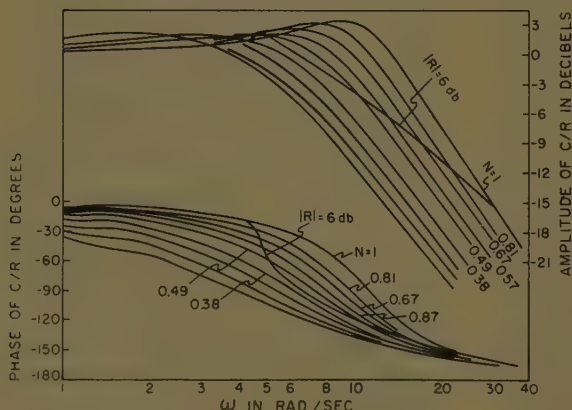


Fig. 8. Replot of Fig. 5 with  $|R| = 6$  db curve drawn

used which might be worth mentioning. From equation 4 it is seen that

$$\frac{X}{R} = \left( \frac{C}{R} \right) \left( \frac{s}{10N} \right) = \left( \frac{C}{R} \right) \left( \frac{j\omega}{10N} \right)$$

Therefore, for any given  $N$  and  $\omega$ , the value  $|X/R|$  in Fig. 6 can be determined from the corresponding value of  $|C/R|$  in Fig. 5 by a simple multiplication (by using logarithms). So Fig. 5 was drawn first using the normalized curves and Fig. 6 was obtained by adding  $20 \log 10N$  to the amplitude curves point by point. This procedure is generally easier as the relation between  $X$  and  $C$  is an open-loop function and usually is much simpler than the closed-loop functions.

To find the amplitude of  $R$  associated with each point on these curves note the identity

$$20 \log |X/R| = 20 \log |X| - 20 \log |R|$$

so that the input  $|R|$  for each curve in Fig. 6 is given by

$$20 \log |R| = 20 \log |X| - 20 \log |X/R|$$

Thus, in order to find the input for a point on the amplitude curves of Fig. 6, take the value of  $|X|$  in decibels corresponding to the  $N$  of the curve in question and subtract from it the value of  $|X/R|$  given by the curve itself at that point.

To illustrate this consider point  $P$  in Fig. 6 for which  $\omega = 5.5$  rad/sec (radians per second),  $|X/R| = 2.7$  db, and  $N = 0.49$ . The value of  $|X|$  corresponding to  $N = 0.49$  is found to be  $X = 8.0^5$  db in Table I; therefore, from equation 7 the amplitude of the input for this point is

$$|R| = 8.0^5 - 2.7 = 5.3^5 \text{ db}$$

It follows that the points on the amplitude and phase curves of Fig. 5 having the same  $N$  and  $\omega$  as point  $P$  must have the same input value and they have been so marked.

In the foregoing example: some arbitrary point on the  $|X/R|$  curves was picked; the values of  $\omega$ ,  $N$ ,  $|X/R|$ , and  $|R|$  for this point were found; and the values were used to determine  $|R|$  for any point. Actually, since contours of constant  $R$  will eventually be wanted, it would be nicer if the required value of  $R$  could be selected and then the points on the  $|X/R|$  curves which corresponded to that value could be identified. These points could be transferred to the  $C/R$  curves as they were drawn and then connected to give a contour for the given input value of  $R$ . This could be done. Since  $N$  is constant along one of the curves of Fig. 6, the value of  $|X|$  is constant along each curve ex-



Table I. Root Locations for  $s^2(s+10)+132N(s+1)=0$

$N$	$\zeta$	$\omega_n$ , sec <sup>-1</sup>	$P_s$	$ X $ , db
2 ... 1	0.4 ... 11.2	1.08 ... $\leq 0$		
5 ... 0.81 <sup>s</sup>	0.45 ... 9.7	1.11 ... 3.0 <sup>s</sup>		
9 ... 0.67 <sup>s</sup>	0.5 ... 8.9	1.13 ... 5.0		
5 ... 0.57	0.55 ... 8.0 <sup>s</sup>	1.16 ... 6.6 <sup>s</sup>		
4 ... 0.49	0.6 ... 7.3	1.19 <sup>s</sup> ... 8.0 <sup>s</sup>		
0 ... 0.38	0.7 ... 6.2	1.29 ... 10.4		
5 ... 0.31 <sup>s</sup>	0.8 ... 5.4	1.43 ... 12.1		
5 <sup>s</sup> ... 0.27	0.9 ... 4.6 <sup>s</sup>	1.63 ... 13.3 <sup>s</sup>		
2 ... 0.24	1.0 ... 1.0	2.00 ... 14.3 <sup>s</sup>		

=1. Let us try to find all the points for which  $|R|$  has the value, say 6 db. Now on the  $N=0.57$  curve, for example, we see from Table I that  $|X|=6.6$ .<sup>5</sup> Therefore, from Equation 6,  $|R|$  will equal 6 at the point on this curve for which

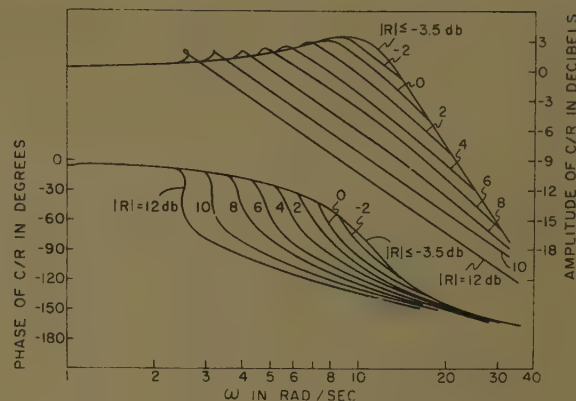
$$|C/R|=6.6^s-6=0.6^s \text{ db}$$

That is, at points where the  $N=0.57$  curve in Fig. 6 crosses the 0.6<sup>s</sup>-db line. There are two such points; see Fig. 6. Similar calculations with other values of  $|R|$  show where other  $|R|=6$  db points occur; these have been identified in Fig.

To prevent any confusion from Fig. 7, it should be pointed out that, although the assumed restrictions on the type of nonlinearity being studied force  $N(X)$  to be a single-valued function of the amplitude of  $X$ , the converse is not generally true. Thus, in the present case,  $N=1$  whenever  $|X|\leq 1$  and the linear  $N=1$  curve is valid whenever  $|X|\leq 0$  db. From equation 6,  $|X/R|$  is given by the  $N=1$  curve at all points for which  $|X/R|\leq |R|$ . This means that a single point on the linear curves may be valid for more than one input amplitude. This should cause no more difficulty than when the system is entirely linear if the nature of the nonlinearity is kept in mind throughout the analysis.

The points on Fig. 7 can be transferred to the curves of Fig. 5, as before, which has been done in Fig. 8, and have been connected to form the contours of constant  $|R|=6$  db. With reference to these curves it can be seen that the input is low enough at low frequencies so that the output is able to follow fairly well. The error,  $E$ , is so small that  $X$  does not saturate. However, as the frequency gets higher with  $R$  constant, the system lags more and more increasing the error until it finally becomes greater than 1, then saturates and  $N$  decreases. In fact, as seen in Fig. 8,  $N$  decreases rather rapidly with frequency in this range. This might be qualitatively explained by the fact that as  $X$  starts to saturate the output has

Fig. 9. Closed-loop frequency response of system of Fig. 2 for various values of input



even more trouble following the input than it would have if the system were linear and the amplitude and phase of the output cannot change much if  $Y$  is able to change only slightly. Therefore, the error increases very rapidly causing  $N$  to decrease rapidly over a short frequency range. (In fact, as will be seen later,  $N$  can change instantaneously in some instances.)

Eventually, the low-pass characteristic of  $G_1$  takes over and  $X$  slowly starts to decrease again. In this region, therefore,  $Y$  is very nearly constant and the output is determined by  $G_2$  (in this case, a straight line with slope of about  $-6$  db per octave). Finally,  $G_1$  makes  $X$  so small that it no longer saturates and operation is again linear.

The curves for other values of input can be obtained in the same way. Fig. 9 shows several such curves drawn for values of  $|R|$  chosen 2 db apart. This system has been simulated on an analog computer and the results are compared with Fig. 9 in Appendix I.

The reader may be interested in the little curlicues that occur in Fig. 9 for high input values. The system was not particularly chosen to show the phenomenon and actually the curve for  $|R|=12$  db is a little unusual in that it shows a closed figure, but this is the so-called jump resonance. If the input amplitude were held at this high value and the frequency increased slowly from a low value, the system would operate at first as previously explained. But as operation continues along the curve, the system reaches a point where the frequency cannot keep increasing with operation staying smoothly on the curve. The output must then jump to the next portion of the curve, with  $X$  increasing and  $N$  decreasing discontinuously, where things again run smoothly. On the other hand, if the frequency were decreased from a high value, and the input amplitude were decreased from a high value and the input amplitude were kept constant at 12 db, a jump would again occur, but at a

different, lower frequency. Thus, for some frequencies there are two possible output amplitudes. (Actually the curve shows three values at these frequencies, but as Levinson<sup>6</sup> points out, the third is unstable.) This phenomenon is discussed by Lozier,<sup>2</sup> Prince,<sup>5</sup> Levinson,<sup>6</sup> and others. Unfortunately, the approximation used by both Lozier and Prince leads to the conclusion that a system, such as the one discussed here, would always exhibit jump resonance for all input amplitudes which is not the case.

Fig. 9 provides all the necessary information for analyzing the sinusoidal behavior of the linearized system of Fig. 2. The thoroughness of the data is limited only by the amount of work the engineer is willing to do. It hardly seems necessary to draw curves closer than 2 db as the depicted curves give a good picture of how the system behaves. Furthermore, in any practical system, it would probably be expected that the input would never exceed some given value, in which case it would be unnecessary to draw curves for higher input values. However, curves can be obtained for as high an input amplitude as desired by drawing enough  $N$  curves.

Finally, in some applications it might be desirable to have normalized output  $|C|$  plotted rather than the output-input ratio. This can be accomplished by simply adding the appropriate value of  $|R|$  to each of the curves of Fig. 9, which has been done in Fig. 10. If the phase position of the input is taken as zero, then the phase of the output is the same as that of  $C/R$  shown in Fig. 9. Since the output of the saturation nonlinearity,  $y$ , approaches a square wave of amplitude 1 as the input,  $x$ , gets large, the curves of Fig. 10 approach a limit which has been included on the graph.

## Systems Containing More Than One Nonlinear Element

Although considerable work has been done and is being done toward finding the

stability, transient response, sinusoidal response, and response to statistically described inputs of systems, which are adequately characterized as having only one significant nonlinear element, little has been attempted on systems in which more than one of the nonlinearities are important.<sup>7</sup> However, such systems do occur, such as a servomechanism in which an amplifier, whose input is limited, drives a motor whose nonlinear characteristics cannot be neglected. This section describes how the previously discussed method can be adapted for the sinusoidal analysis of systems of this sort.

It is not hard to see why much study has not been done on systems containing several nonlinearities. In the first place, quite often all but one of the nonlinearities are negligible, even though this is not always the best design. But more important, it is hard enough to work with only one nonlinear element and adding a second not only doubles the work but compounds the difficulties.

This paper utilizes the describing function to characterize the nonlinear elements. The input to each nonlinearity is assumed to be a sinusoid when the system input is a sine wave. With the single nonlinearity it is fairly common for the open-loop characteristic to be low-pass because of mechanical devices in the loop, slow amplifiers, etc., but this is not sufficient when there are several nonlinear elements. If describing functions are used, there must be a low-pass filter between every two nonlinearities so that the harmonics created by one are sufficiently filtered to make the input to the next essentially sinusoidal. Thus, a describing function method generally will be less valid in the multiple nonlinearity case and the accuracy may leave something to be desired. However, if there is at least one integration,  $(1/s)$ , between the nonlinearities, the method should provide at

least an indication of the sinusoidal response characteristics of the system. (Actually, although the system given in the following example was picked completely at random, the results of an analog computer study are quite close to the analytical calculation.)

A simple system containing two nonlinear elements is shown in Fig. 11 where the linear elements are expressed in terms of their Laplace transfer functions,  $G_1$ ,  $G_2$ ,  $G_3$ , and  $H$ , and the nonlinear elements are given by their describing functions  $N_1$ ,  $N_2$ . Assuming that the describing functions are valid, it is seen that

$$\frac{C}{R} = \frac{G_1 N_1 G_2 N_2 G_3}{1 + G_1 N_1 G_2 N_2 G_3 H} \quad (8)$$

$$\frac{X_1}{R} = \frac{G_1}{1 + G_1 N_1 G_2 N_2 G_3 H} \quad (9)$$

With only one nonlinearity the poles and zeros of the closed-loop functions were found as a function of the gain  $N$  and curves of  $C/R$  and  $X/R$  were plotted. Since  $N$  and  $X$  had a known functional relationship (the describing function), it was possible to find the value of  $R$  corresponding to each point on the  $X/R$  curves. These points could then be identified on the  $C/R$  curves which provide the final curves of  $C/R$  for constant  $R$ . In the present case, the poles and zeros may still be found as a function of gain, but the gain is now the product of the two describing functions,  $N_1 N_2$ , and is not simply related to any of the signal amplitudes.

Therefore, to find the input value for each point, it is necessary to know how  $N_1 N_2$  is related to the amplitude of  $X_1$ . This may be done in two ways. The more obvious way is to choose several values of  $X_1$  first. For any given amplitude of  $X_1$ , the amplitudes of  $N_1$  and  $Y_1$  are independent of the frequency and are obtainable from the describing function  $N_1$

( $X_1$ ). For any given frequency of  $Y_1$  the amplitude ratio,  $|X_2/Y_1|$ , is independent of amplitude and is known from the transfer function  $G_2$ . Thus, for any given amplitude and frequency of  $X_1$ , it is possible to find the corresponding amplitude of  $X_2$  and therefore  $N_2$ . For the given  $X_1$  the product,  $N_1 N_2$ , is then known. If this process is carried through for a large number of values of  $X_1$ , curves could be drawn from which  $X_1$  could be determined by any given  $N_1 N_2$  and  $\omega$ .

This relationship may also be obtained by choosing  $X_2$  first. For any given  $X_2$  the value of  $N_2$  is independent of frequency and can be determined from the describing function  $N_2(X_2)$ . Thus, for any given amplitude and frequency of  $X_2$ , it is possible to find the corresponding amplitude of  $Y_1$ , by using  $G_2$  inversely. Now, though it is not usually thought of in this way, the describing function is mathematically just as much a function of an element's output as its input. Certainly, under the assumption that the input is a sinusoid, if the fundamental of the output is given, the value of the describing function can be found. Therefore, for the given amplitude and frequency of  $X_2$ , if  $Y_1$  is known, then both  $N_1$  and  $X_1$  can be determined and the corresponding value of the product,  $N_1 N_2$ , can be calculated. Again, if this is done for several values of  $X_2$  and  $\omega$ , curves could be drawn relating  $X_1$ ,  $N_1 N_2$ , and  $\omega$ . This second approach is used here because it yields data in a form that is a little easier to handle; thus, more accuracy is possible with less effort.

There appears to be no obvious reason why  $X_1/R$  curves had to be used to find the necessary relations for the determination of the desired input-output curves. With a single nonlinear element it was very convenient to do so since  $N$  was constant along any one of the curves. However, in the present case this is not true. From the previous discussion it is fairly obvious that it would not be as easy to relate  $N_1 N_2$ ,  $\omega$ , and  $X_1$  then to relate  $N_1 N_2$ ,  $\omega$ , and  $X_2$ . So,  $Y_2/R$  curves could be used instead of the  $X_1/R$

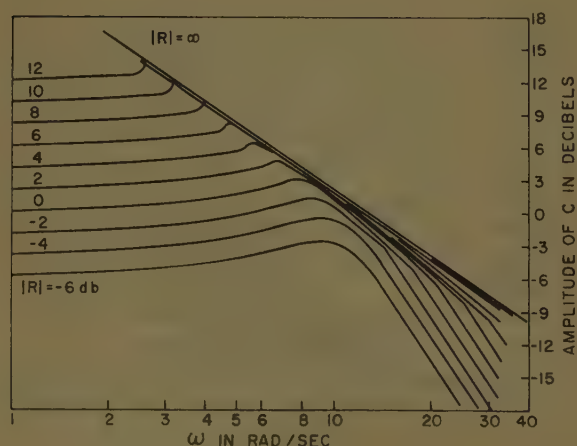


Fig. 10. Output of system of Fig. 2 for various values of input

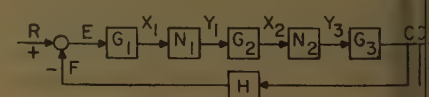


Fig. 11. System containing two nonlinear elements

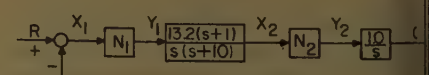


Fig. 12. Specific example of system containing two nonlinear elements



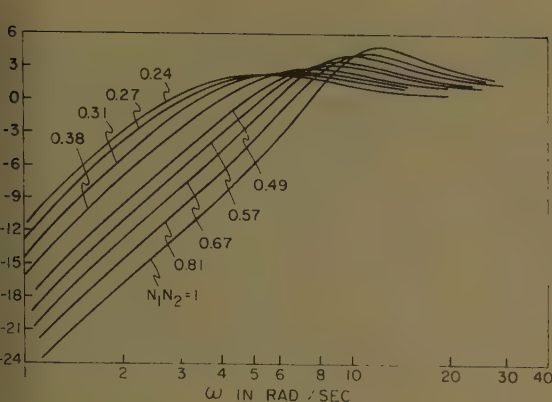


Fig. 13.  $X_1/R$  curves for system of Fig. 12

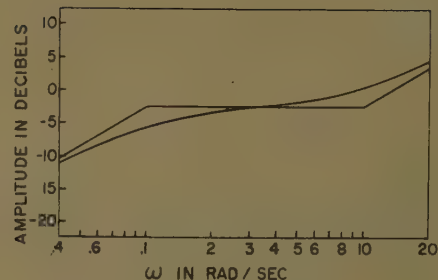


Fig. 14. Amplitude of  $G_2^{-1}$

$$= \frac{j\omega(j\omega+20)}{13.2(j\omega+1)}$$

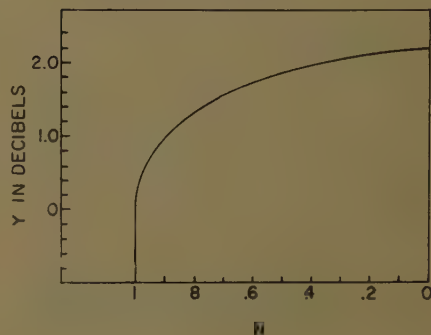


Fig. 15. Inverse describing function for saturation

curves, and a little more effort would reveal the relation between  $N_1N_2$ ,  $\omega$ , and itself, in which case the original  $C/R$  curves could be used without any recourse to a second set such as those of equation 9. Actually the amount of work needed in any of these cases is just about the same; the choice of method based on greater convenience and other such considerations. In the following example,  $X_1$  is used rather than  $Y_2$  or  $C$  as it is somewhat easier and more accurate.

To illustrate how the input-output curves are obtained, consider the system

shown in Fig. 12 where both  $N_1$  and  $N_2$  are taken to be saturation as shown in Fig. 2(B). The gain values have been so chosen that when operating in the linear region on  $N_1$ , ( $X_1 \leq 1$ ), the system is the same as that studied in the previous section. Therefore, the curves of Fig. 5 are the  $C/R$  curves for this system if  $N$  is replaced by  $N_1N_2$ . The curves of the amplitude of  $X_1/R$  are shown in Fig. 13.

Now, since  $X_2$  will be chosen and then  $Y_1$  determined, a plot of

$$\left| \frac{Y_1}{X_2} \right| = \left| \frac{1}{G_2(j\omega)} \right| = |G_2^{-1}(j\omega)| = \left| j\omega \frac{(j\omega+10)}{13.2(j\omega+1)} \right| \quad (10)$$

will be needed. Because  $N_1$  and  $N_2$  are only amplitude sensitive, just the amplitude of  $G_2^{-1}$  is required and this is drawn in Fig. 14. Further, a curve relating  $N_1$  to  $Y_1$  is needed; this is just a redrawing of Fig. 3 using the fact that

$$Y_1 = N_1(X_1)X_1 \quad (11)$$

and is shown in Fig. 15.

Now the values of  $|X_2|$  must be selected. The values of the product,  $N_1N_2$ , that are wanted are the values of  $N$  given in Table I. Since  $N_1=1$  for  $|X_1| \leq 1$ , let values of  $|X_2|$  be chosen such that  $N_2$  takes on the values of  $N$  given in Table I. Then  $N_1N_2$  will have these values when  $N_1$  is not saturating. Also, select a few more values between these just to increase the number of points on the curves. Table II provides the values of  $|X_2|$  and  $N_2$  that will be used in this discussion.

Table II. Describing Function for Saturation from Fig. 3

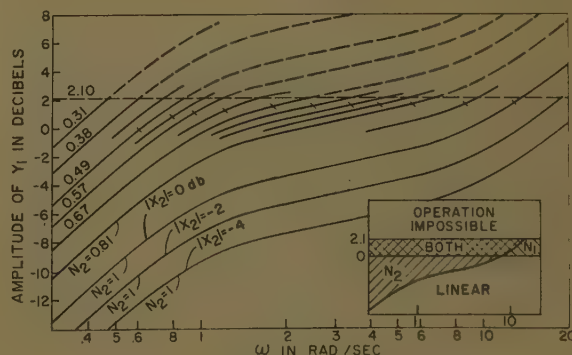
$X_{21}$ , db	$N_2$
0	1
1.85	0.90
3.05	0.815
4.55	0.75
5.85	0.75
7.35	0.71
8.05	0.675
9.1	0.62
10.4	0.57
12.1	0.53
13.35	0.49
14.35	0.44
15.85	0.38
17.35	0.315
18.85	0.27
20.4	0.24

Table III. Points on Fig. 16 Where  $N_1N_2=0.49$

$N_2$	$N_1$	$ Y_1 $ , db	$\omega$ , rad/sec	$ X_1 $ , db
0	0.49	1.86	13.1	8.0
0.05	0.545	1.78	8.7	7.1
0.1	0.60	1.72	5.68	6.2
0.15	0.63	1.68	4.36	5.7
0.2	0.65	1.64	3.42	5.3
0.25	0.69	1.58	2.52	4.8
0.3	0.725	1.51	1.81	4.3
0.35	0.79	1.35	1.23	3.4
0.4	0.86	1.14	0.95	2.45
0.45	0.925	0.83	0.90	1.55
0.5	1	$\leq 0$	$\leq 0.61$	$\leq 0$

impossible points for  $N_1N_2=0.49$ ;  
 $N_2 < N_1N_2$ , since  $N_1 \leq 1$ .

Fig. 16. Amplitude of  $Y_1$  for system of Fig. 12. Insert shows where the two nonlinear elements saturate



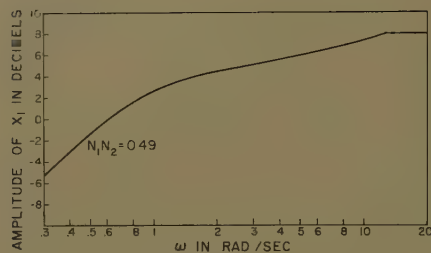


Fig. 17. Amplitude of  $X_1$  for system of Fig. 12 and  $N_1N_2=0.49$

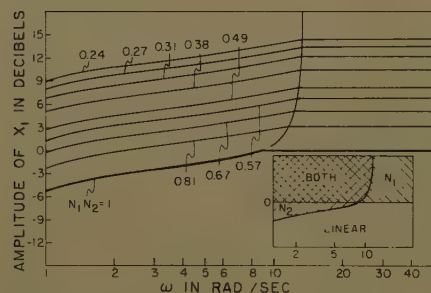


Fig. 18. Amplitude of  $X_1$  for system of Fig. 12. Insert shows where the two nonlinear elements saturate

Thus, the regions of Figure 16 in which each nonlinear element is saturating may be determined.

Now the points in Fig. 16 where the product  $N_1N_2$  has the values used in Figs. 5 and 13 may be found. In the region where only  $N_2$  saturates,  $N_1=1$ ,  $N_1N_2=N_2$ , and the desired  $N_1N_2$  curves are those already drawn in Fig. 16. In the region where only  $N_1$  saturates, the product,  $N_1N_2$ , is entirely determined by  $N_1$  which in turn is a function of  $|Y_1|$ . In this region the constant  $N_1N_2$  curves are constant  $|Y_1|$  curves which are horizontal lines.

In the region where both nonlinearities saturate, some point-by-point calculation will be needed to find the desired points. For instance, let us find the points where  $N_1N_2=0.49$  by using Table III. The first column gives the values of  $N_2$  on the various curves of Fig. 16 as given in Table II. The second column is merely  $N_1N_2=0.49$  divided by  $N_2$ . The

third column is obtained from the inverse describing function of Fig. 15 and gives the value that  $|Y_1|$  must have to yield the  $N_1$  given in the previous column. The first three columns contain all the necessary data to find the  $N_1N_2=0.49$  points on Fig. 16.  $N_1N_2$  must equal 0.49 at each point where the  $N_2$  curve in column 1 intercepts the value of  $|Y_1|$  given in column 3. These points have been identified in Fig. 16.

Now the curves relating  $X_1$ ,  $\omega$ , and  $N_1N_2$  are needed, so columns 4 and 5 have been added to Table III. The frequencies in the fourth column are read directly from Fig. 16 and are the frequencies at which the various points occur. The fifth column is obtained from the describing function of Fig. 3. Now the curve of  $|X_1|$  versus  $\omega$  may be drawn for  $N_1N_2=0.49$  and is given in Fig. 17. For  $\omega \leq 0.61$  rad/sec this curve is identical with the  $N_2=0.49$  curve of Fig. 16 because  $X_1=Y_1$ . For  $0.61 \leq \omega \leq 13.1$  rad/sec the curve is plotted from columns 4 and 5 of Table III, and for  $\omega \geq 13.1$  rad/sec  $N_1=0.49$  and  $|X_1|=8.0$  db.

This curve is repeated in Fig. 18 together with similar curves obtained in the same manner for the other values of  $N_1N_2$ . Fig. 18 is drawn to the same scale as the  $|X/R|$  curves of Fig. 13. It is important for subsequent development that these two curve sheets be drawn to the same scale.

Now it is possible to find the value of the input  $R$  at any point on the curves of Fig. 13. These curves give  $|X_1/R|$  versus  $\omega$  for a certain set of values of gain  $N_1N_2$ , while Fig. 18 gives curves of  $|X_1|$  versus  $\omega$  for the same values of  $N_1N_2$ . Therefore, from the identity

$$20 \log |R| = 20 \log |X_1| - 20 \log |X_1/R| \quad (12)$$

the value of the input,  $|R|$ , for any given  $N_1N_2$  and  $\omega$  is the value of  $|X_1|$  in db, given by Fig. 15, minus the value of  $|X_1/R|$  in db, given by Fig. 13. There are several ways to perform this subtraction. It is certainly valid just to pick points at random and use equation 12 to find  $|R|$

at these points, but as with a single nonlinearity, it is better and possible to choose the value of input first and then find the appropriate points on the curves. The first thing that comes to mind is some sort of trial-and-error method where a few points are selected with an eye toward converging on the desired point by the application of equation 12, but a more systematic method is available if the curves of  $|X_1/R|$  and  $|X_1|$  are drawn on tracing paper using the same scale. By use of equation 12, for example, the 6 db points occur where the  $|X_1/R|$  curve and the  $|X_1|$  curve for the same  $N_1N_2$  are 6 db apart, with the latter greater. So, one set of curves is placed over the other with the 0-db lines coincident and a pair of dividers is set at 6 db, these points can be found with very little difficulty.

Actually there is an even better way of doing this same thing. If the  $|X_1/R|$  curves are placed over the  $|X_1|$  curves so that the 0-db line of the  $|X_1/R|$  curves coincides with the 6-db line of the  $|X_1|$  curves, then equation 12 is automatically satisfied at intersections of similar  $N_1N_2$  curves for  $|R|=6$  db. An attempt to show this process is given in Fig. 19 where the resulting 6-db points are shown. These points may then be transferred to the  $C/C$  curves and connected to form the  $|R|=6$ -db contour, which has been done in Fig. 20.

Similar statements may be made about the  $R=6$ -db curves of Fig. 20 as were made about Fig. 8. If the amplitude of the input is kept constant at  $R=6$  db and the frequency is slowly increased from a low value, then at the low frequencies, the error is small and the system operates linearly along the  $N_1N_2=1$  curve. Eventually, the error increases so that one of the nonlinearities begins to saturate. In this particular instance, the second nonlinearity,  $N_2$ , saturates first. Now, with only  $N_2$  saturating, this system is the same as that with only one nonlinearity, and this portion of the curve is identical with that of Fig. 8. As before, the error increases and the gain decreases rapidly because  $Y_2$  cannot increase rapidly. Now, the first nonlinear element,  $N_1$ , begins to saturate, but it does so in such a way that the system becomes unstable immediately. The error must increase instantaneously and the gain decreases instantaneously with a slight increase of frequency so that the output jumps to the upper section of the curve. This is the same sort of effect as the jump response mentioned previously but the effect caused here by the presence of two saturations.

With further increase of frequency,

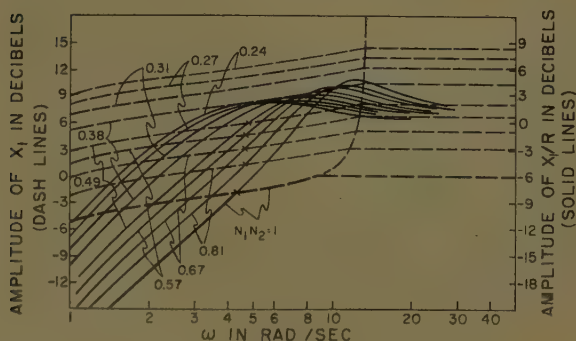


Fig. 19. Determination of  $|R|=6$  db points



pass characteristics of both the  $G_2 = 2(s+1)/s(s+10)$  and  $G_3 = 10/s$  begin to take effect and both  $X_2$  and  $C$  decrease regardless of the fact that the error,  $X_1$ , remains large. Eventually,  $X_2$  becomes small that it no longer saturates and only  $X_1$  remains saturated. The point at which this occurs lies on the curve separating the regions where both  $N_1$  and  $N_2$  saturate and only  $N_1$  saturates. This curve was taken directly from Fig. 18 and is shown in Fig. 20. At higher frequencies, the output  $C$  continues to decrease and  $X_1$  approaches the input  $R$ . Thus, the first nonlinearity,  $N_1$ , remains saturated at all high frequencies for  $|R| = 6$  db.

In a similar manner, curves like those of Fig. 20 can be obtained for other input values. Fig. 21 shows several of these curves. This system has been simulated on an analog computer and the results are compared in Fig. 21 in the Appendix. Since this section has dealt with graphical solution it is well to discuss briefly the accuracy of the results. At several points in the work, pairs of curves have been graphically added to obtain new curves. Thus, the accuracy obtained by such techniques is questionable and scepticism is not without ground; however, the real question is how much accuracy is warranted. The method is based upon the assumption that the inputs to the various nonlinear elements are all sinusoidal when the input is sinusoidal. This is only approximately true even in the case of a single nonlinearity and when there are several nonlinear elements, the approximation is much worse. In most instances, the method will, at best, probably give only an indication of how the output varies with the amplitude and frequency of the input. The curves should be drawn as accurately as possible under the assumptions, but it hardly seems necessary to require more accuracy than can be obtained from graphical techniques.

## Conclusions

A method has been developed for finding the response of certain types of nonlinear feedback control systems to steady-state sinusoidal input signals. The method makes use of the describing function to represent the nonlinear elements. Familiar linear frequency response techniques are used to draw families of curves from which the input-output data of the nonlinear system are deduced. For any given system the method is no more difficult than the plotting of the frequency response curves of the same system with the nonlinearities replaced by linear gain

elements for several values of gain. The nonlinear response data are obtained in the form of contours of fundamental (and any desired harmonics; see Appendix II) of the output signal as a function of frequency and the amplitude of the system input signal. Complete response data are available for any desired frequency or input amplitude from a single set of curve sheets. A new set of curves is not needed for every different input amplitude.

The method is illustrated by means of two examples which are worked out in full. As a check on the accuracy of the results the two systems have been simulated on an analog computer. The data from the computer roughly show the method to be as accurate as the computer itself. The second example involved a system containing two nonlinear elements and, although the method would not be expected to give as accurate results for this case as for the single nonlinearity case, the analog results agree very closely with the curves obtained analytically.

## Appendix I. Analog Computer Simulation

The systems used as examples have been simulated on an analog computer. The results of the simulation are shown in Figs. 22 and 23.

Fig. 20. Fig. 5 replotted with  $|R| = 6$  db curve drawn. Curves A and B show where  $N_2$  just begins to saturate

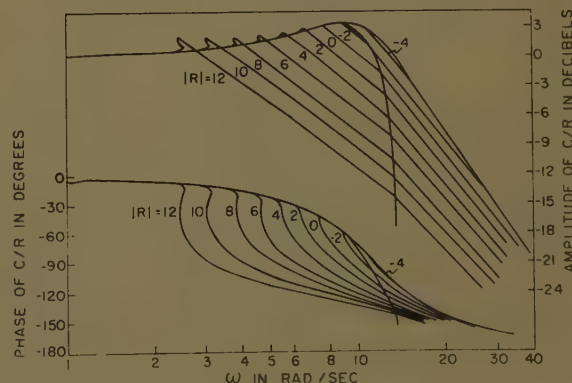
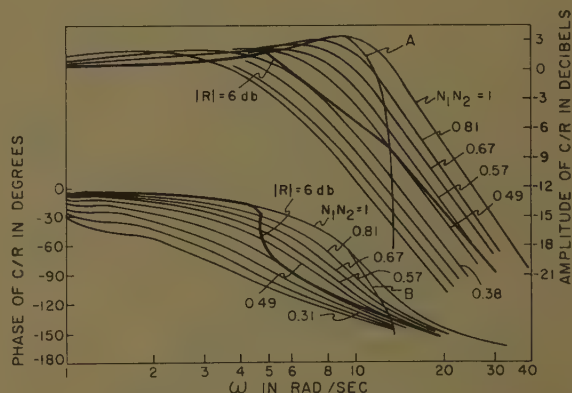


Fig. 21. Closed-loop frequency response of system of Fig. 12 for various values of input

## Appendix II. Harmonics in the Output

As was stated previously, although the input to the nonlinear element may be so nearly sinusoidal that the describing function may be used, the system output may contain significant harmonics and it may be desirable to calculate their values. This is not at all difficult to do when the describing function is indeed applicable. Truxall<sup>1</sup> has shown that for idealized saturation if the input is a sinusoid of amplitude  $|X|$  the  $n$ th harmonic of the output is

$$|Y^{(n)}| = \frac{2|X|}{\pi n} \left\{ \frac{\sin(n-1)t_2}{n-1} + \frac{\sin(n-1)t_2}{n-1} \right\} \quad (13)$$

for  $|X| > 1$  and  $n = 1, 3, 5, \dots$ , where  $t_2 = \arcsin[1/|X|]$  and the notation  $Y^{(n)}$  is taken to mean the  $n$ th harmonic of the output  $y$ . Thus the harmonics of  $y$  are known and since the transfer function  $G_3$  from  $y$  to the system output  $c$  is linear the harmonics in the output can easily be calculated.

To illustrate how this is done let us consider the third harmonic output of the system considered in the preceding section.

$$\left| \frac{Y'''}{X} \right| = \frac{1}{3\pi} \left( \sin 2t_2 + \frac{1}{2} \sin 4t_2 \right) \quad (14)$$

Now Table IV can be filled in. Column 1 contains the values of  $N$  used throughout this analysis; column 2 gives the corresponding values of  $|X|$  from Table I; column 3 provides the corresponding values of  $|Y'''|/|X|$  from equation 13; and column 4 contains the sum of columns 2 and 3. For any given input frequency,  $\omega$ , therefore, the

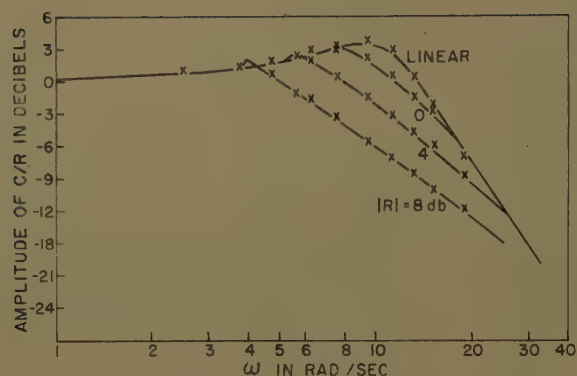


Fig. 22. Results of analog simulation of system of Fig. 2

Solid lines repeated from Fig. 9  
x represents analog data

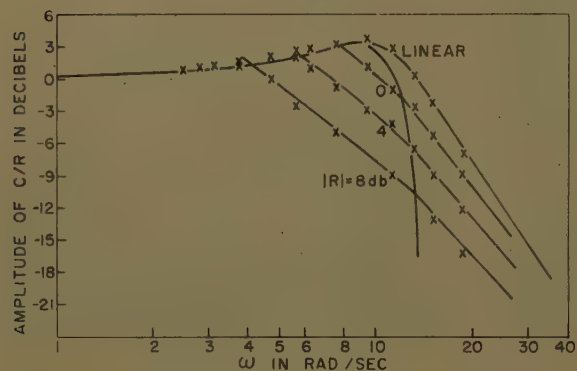


Fig. 23. Results of analog simulation of system of Fig. 12

Solid lines repeated from Fig. 21  
x represents analog data

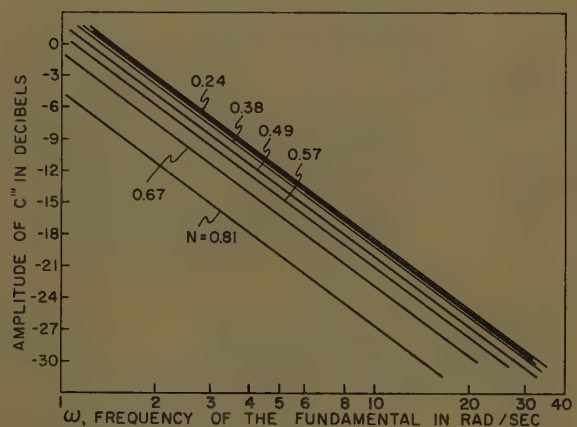


Fig. 24. Third harmonic output of system of Fig. 2 for various values of gain N

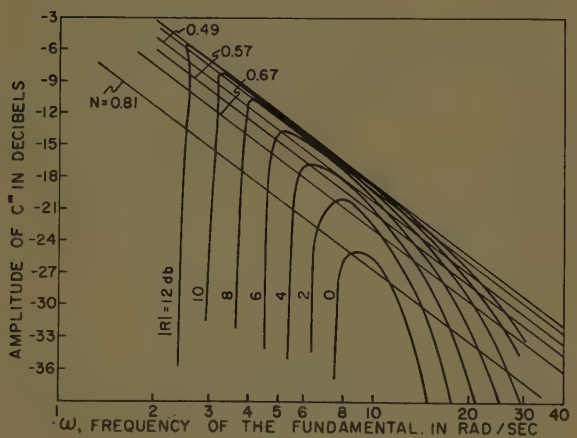


Fig. 25. Third harmonic output for various values of system input

third harmonic component of system output (see Fig. 2) is

$$C''' = G_2(j3\omega)Y''' = \frac{10}{j3\omega} Y''' \quad (15)$$

The curve of  $|G_2(j3\omega)| = 10/j3\omega$  is a

straight line of slope minus 20 db, per decade. Therefore, for constant values of  $|Y'''|$  corresponding to constant values of  $N$ , equation 15 is a series of straight lines as shown in Fig. 24 where the abscissa is the frequency of the system input, that is, the frequency of the fundamental of the output.

Table IV. Third Harmonic in Saturation

$N$	$ X' $ , db	$ Y''' $ , db	$ Y $ , db
1	$\leq 0$	$-\infty$	$-\infty$
0.81 <sup>a</sup>	3.0 <sup>a</sup>	-19.2	-16.1
0.67 <sup>a</sup>	5.0	-17.3 <sup>a</sup>	-12.3
0.57	6.6 <sup>a</sup>	-17.3 <sup>a</sup>	-10.0
0.40	8.0 <sup>a</sup>	-17.7	-9.0
0.38	10.4	-19.0 <sup>a</sup>	-8.6
0.31 <sup>a</sup>	12.1	-20.4	-8.3
0.27	13.3 <sup>a</sup>	-21.5	-8.1
0.24	14.3 <sup>a</sup>	-22.4	-8.0

Now, from the work in the paper, the value of  $N$  and  $\omega$ , corresponding to given values of input  $R$ , are known and can be identified in Fig. 24. This has been done on Fig. 25 and the points have been connected to form contours of constant input.

The relative phase of third harmonic requires discussion. From equation 15 it is seen that  $C'''$  lags  $Y'''$  by  $\pi/2$  radians. This means that when the third harmonic of  $y$  is going through zero with positive slope, the third harmonic of  $c$  is going through its negative maximum. To see how this relates to the phase of the fundamental note that the output waveshape of the nonlinear element has sine symmetry (is an odd function) about the point where it passes through zero. Therefore, when the fundamental of  $y$  is going through zero with positive slope, all the harmonics also pass through zero and with positive slope. Furthermore, at this same instant the fundamental of the output,  $c$ , is going through its negative maximum as are all the harmonics in the output. Thus, the output waveshape has cosine symmetry (is an even function) about this point and the output wave is going through its negative maximum. This means that the output waveshape will always be peaked, that is, more pointed than a sine wave.

Comparing Fig. 25 with Fig. 10 it is seen that the amplitude of the third harmonic of the output is more than 19 db lower than that of the fundamental. That is, the third harmonic output has about 1/10 the amplitude of the fundamental or smaller. Furthermore, the higher harmonics would be even less significant. The output waveshape is itself not very different from being sinusoidal, which was assumed.

In addition, the factor  $1/s$  in  $G_1$  would decrease the harmonics still further so that the input to the nonlinearity is indeed very nearly sinusoidal. This harmonic analysis, as the analog study of Appendix I confirms, shows that neglecting harmonics is permissible and the method presented should give highly accurate results.

## References

1. AUTOMATIC FEEDBACK CONTROL SYSTEMS, THESIS (book), J. G. Truxall. McGraw-Hill Book Company, Inc., New York, N. Y., 1955.
2. A STEADY STATE APPROACH TO THE THEORY OF SATURABLE SERVO SYSTEMS, *Proceedings, Institute of Radio Engineers*, New York, N. Y., 1956, PGAC-1, May 1956, pp. 19-39.
3. AN ANALYTIC METHOD FOR FINDING THE CLOSED-LOOP FREQUENCY RESPONSE OF NONLINEAR FEEDBACK-CONTROL SYSTEMS, K. Ogata.



EE Transactions, pt. III (Power Apparatus and Systems), vol. 76, Nov. 1957, pp. 277-85.

A GENERAL METHOD FOR ANALYZING AND SYNTHESIZING THE CLOSED LOOP RESPONSE OF A LINEAR AND A NONLINEAR SERVOMECHANISM, H. El-Sabbagh. WESCON Convention Record, Institute of Radio Engineers, pt. 4 (Automatic Control), 1957, pp. 58-77.

GENERALIZED METHOD FOR DETERMINING THE CLOSED-LOOP FREQUENCY RESPONSE OF NONLINEAR SYSTEMS, L. T. Prince, Jr. AIEE Transactions, pt. II (Applications and Industry), vol. 73, pt. 1954, pp. 217-23.

SOME SATURATION PHENOMENA IN SERVOMECHANISMS WITH EMPHASIS ON THE TACHOMETER

STABILIZED SYSTEM, E. Levinson. Ibid., vol. 72, Mar. 1953, pp. 1-9.

7. FREQUENCY RESPONSE OF NONLINEAR CLOSED-LOOP FEEDBACK CONTROL SYSTEMS, S. L. Mikhail, G. H. Fett. Ibid., vol. 77, Nov. 1958, pp. 436-38.

8. PRINCIPLES OF SERVOMECHANISMS (book), G. S. Brown, D. P. Campbell. John Wiley & Sons, Inc., New York, N. Y., 1948.

9. SERVOMECHANISM ANALYSIS (book), G. J. Thaler, R. G. Brown. McGraw-Hill Book Company, Inc., New York, N. Y., 1953.

10. SINUSOIDAL ANALYSIS OF FEEDBACK-CONTROL SYSTEMS CONTAINING NONLINEAR ELEMENTS, E. Calvin Johnson. AIEE Transactions, pt. II

(Applications and Industry), vol. 71, Jan. 1952, pp. 169-81.

11. AN EXTENSION OF THE ROOT LOCUS METHOD TO OBTAIN CLOSED-LOOP FREQUENCY RESPONSE OF FEEDBACK CONTROL SYSTEMS, A. S. Jackson. Ibid., vol. 73, Sept. 1954, pp. 176-79.

12. PREDICTION OF TRANSIENT RESPONSE OF NONLINEAR SERVOMECHANISM BY SINUSOIDAL ANALYSIS, H. C. Brearley, Jr. Ph.D. Thesis, The University of Illinois, Urbana, Ill., 1954.

13. A FREQUENCY RESPONSE METHOD FOR ANALYZING AND SYNTHESIZING CONTACTOR SERVOMECHANISMS, Ralph J. Kochenburger. AIEE Transactions, vol. 69, pt. I, 1950, pp. 270-84.

## Discussion

Sridhar (Purdue University, Lafayette, Ind.): It seems quite obvious that most authors who use the describing function of a nonlinearity to obtain the closed loop frequency response of a nonlinear system do not seem to appreciate the full limitations of the describing function. A sufficient condition for the describing function method of obtaining the closed loop frequency response of a nonlinear system to be valid is that the system be totally stable. Otherwise, it is possible to obtain quite erroneous results using this method. It is possible for a system which is perfectly stable for no inputs, as indicated by the describing function method, to break into auto-oscillations when subjected to sinusoidal inputs, the frequency of the auto-oscillations being independent of the forcing frequency. Again, it is possible to quench the auto-oscillations of certain autonomous nonlinear systems by subjecting them to sinusoidal inputs so that the system will respond to the forcing frequency. This latter is, of course, the well-known phenomenon of signal stabilization. It is high time that authors stopped proposing new methods of obtaining the fre-

quency response of nonlinear systems using the describing function method. It will possibly be worth while to establish criterion which will tell when the describing function method is valid and when it is not.

The author appears to have naively used the describing function to obtain the frequency response of a nonlinear system with two nonlinearities. Actually, the validity of the describing function method to investigate the stability of such systems, even when autonomous, has not been established for a general case. I shall appreciate any comments the author has regarding the assumptions on the linear elements when he analyzed the system.

A. S. McAllister: I want to thank Mr. Sridhar for his thoughtful discussion. The points made by him are in most cases well taken. The blind use of linearization techniques can indeed lead to erroneous results. But this has long been the case whether using Hooke's law, studying Class A amplifiers or analyzing feedback control systems. The engineer should always keep in mind that he is making an approximation and should check his results and assumptions frequently and as carefully as possible. If a rigorous

validity criterion that can be applied easily is available it would be well to use it.

All this certainly limits the general applicability of such techniques. However, it does not invalidate the usefulness of linearization. Indeed, linearization is usually easier and faster to use than the more elaborate and rigorous methods. There are times when the linearizing method will give an approximate answer where the exact method becomes so involved as to be utterly useless. The more rigorous methods simply have not yet been developed to the point where they can be used in the more complicated problems.

As for the assumptions made in analyzing a two-nonlinearity system, they are given in the third paragraph of the appropriate section of the paper and need not be repeated here. In this connection, it might be noted that Mikhail and Fett<sup>1</sup> also analyzed a two-nonlinearity system using the two describing functions. They obtained rather good results even though the system they studied had no linear elements, low-pass or otherwise, between the nonlinear elements.

## REFERENCE

1. See reference 13 of the paper.

# On Stabilization of Feedback Systems Affected by Hysteresis Nonlinearities

A. K. MAHALANABIS  
NONMEMBER AIEE

PRESENCE OF HYSTERESIS or hereditary types of nonlinearities leads to small-signal destabilization in servomechanisms and other feedback control systems and produces small-amplitude sustained oscillations or chattering, particularly when precision requirements make necessary the use of a high loop gain. Some safeguards against this undesirable characteristic are therefore presented here.

Several methods of providing proper compensation for hysteresis effects have been suggested.<sup>1-5</sup> The use of coulomb

friction for this purpose has been investigated by this author,<sup>6,7</sup> and the present paper embodies results of further studies in this direction. Some nonlinear compensation schemes for improving the small-signal stability of feedback control systems are dealt with. They not only permit compensation for the action of hysteresis effects on stability but also result in improved response characteristics.

The general desirability of such nonlinear compensation methods in inherently nonlinear systems has been stressed, for

example, by Truxal, who considered an interesting case of nonlinear compensation, around a feedback loop, for backlash effects.<sup>8</sup>

## Hysteresis Effects in Simple Servo Systems

Simple electromechanical systems are considered here. Basic components of an electromechanical servomechanism, shown in Fig. 1, are: (1) the error-sensing unit, (2) the controller-amplifier unit, (3)

Paper 61-709, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE-AIEChE-ASME-IRE-ISA Joint Automatic Control Conference, Boulder, Colo., June 28-30, 1961. Manuscript submitted September 8, 1960; made available for printing May 8, 1961.

A. K. MAHALANABIS is with the University College of Technology, Calcutta, India.

The author is grateful to Professor J. N. Bhar and Dr. A. K. Choudhury for their assistance; also to the Ministry of Scientific Research and Cultural Affairs, Government of India, for providing funds during the progress of this work.

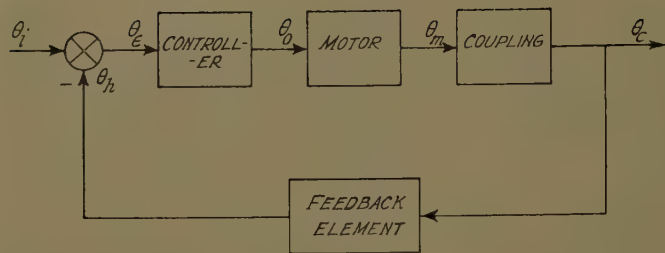


Fig. 1. Basic components of a servo-mechanism

the motor, (4) the output coupling unit.

One or more of these units generally have a hysteresis-type nonlinearity, whose source could be: (1) backlash in the output coupling unit, (2) the controller in a contactor system which incorporates a contactor, e.g., a biopolarized (electromagnetic) relay, possessing hysteresis, and (3) field-controlled electromagnetic devices, e.g., amplidyne generators in the controller unit or a field-controlled motor. That hysteresis may arise from sources 2 or 3 is evident. Backlash in the output coupling produces a hereditary effect when the load has negligible inertia and damping, equal to or greater than critical.

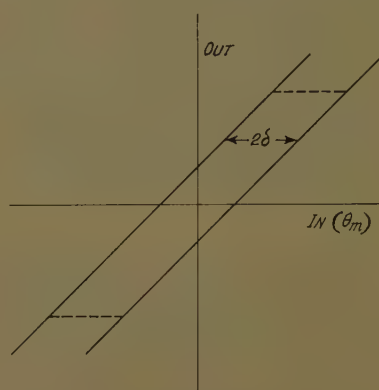
The functional characteristics arising from nonlinearity in different cases are illustrated in Figs. 2(A) through 2(D). Fig. 2(E) is an approximation of the magnetic hysteresis curves, which simplifies analysis by describing function techniques. Performance of the system shown in Fig. 1 is difficult to analyze when more than one form of hysteresis occurs simultaneously. Analyses of individual effects of these nonlinearities have, however, been made. The frequency-response approach, in which the nonlinearity is replaced by its describing function, has been found most convenient for stability analysis of nonlinear systems. The describing functions, Figs. 2(A), (B), and (E), have been discussed in published literature.<sup>1,2,9</sup> Although these functions differ in detailed properties, a common feature is the predominance of nonlinear effects at comparatively small-signal amplitudes, where a reduction in the magnitude of transfer gain is produced, together with a lagging phase shift. Both effects increase as the signal amplitude approaches the hysteresis width.

The phase lags cause the system's stability to be impaired at small-signal levels, producing sustained oscillations of essentially small amplitude in an otherwise stable system.

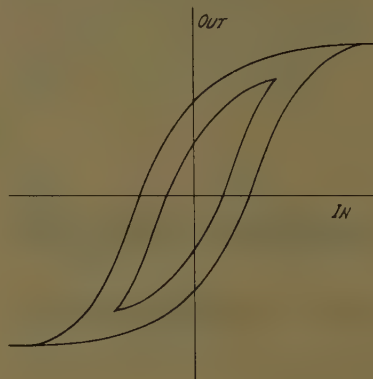
### Stabilization of Oscillations Arising from Hysteresis Effects

Sustained oscillations outlined in the previous section can be prevented. This

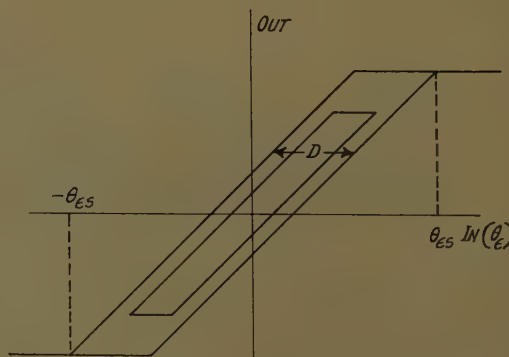
calls for reducing the loop-gain sufficiently to avoid an intersection in the gain-phase shift plane of the linear-frequency locus and the amplitude locus of the negative reciprocal of the relevant describing function. Sluggish response characteristics result from this procedure, however. The nonlinear-gain element suggested



(A)



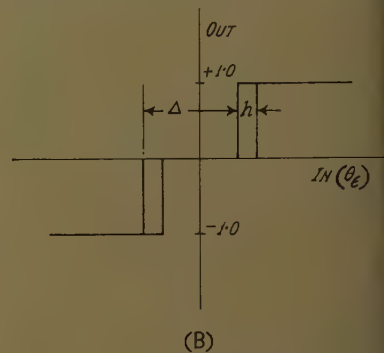
(C)



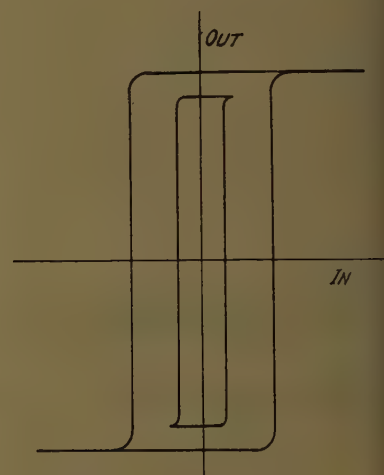
(E)

by Johnson<sup>2</sup> seeks to remedy this defect (see Fig. 3(A)). But an extreme use of this plan, as exemplified by the dead-zone element shown in Fig. 3(B), may impair the static accuracy of the resulting system. Use of conventional phase-lead networks may also be somewhat unsatisfactory, as will be discussed subsequently.

Stabilization of oscillations in systems having hysteresis may be accomplished through improving system stability at small-signal levels without affecting the response characteristics at large-signal amplitudes. This is because hysteresis introduces a sort of variable damping, with effective damping becoming smaller as the signal amplitude diminishes. A logical procedure for compensating for the



(B)



(D)

Fig. 2. Type of hysteresis. A—Caused by backlash. B—Caused by contactor. C and D—Caused by electromagnetic devices. E—Approximation of magnetic hysteresis



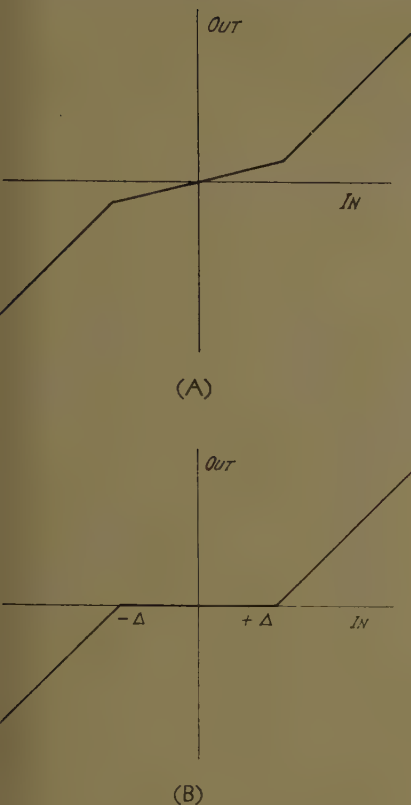
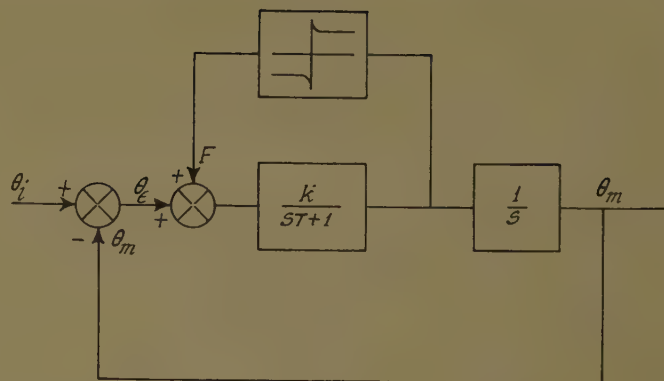


Fig. 3. A—Nonlinear-gain characteristics. B—Dead-zone characteristics

hysteresis effects is, therefore, to insert a component in the forward or feedback path that, in absence of hysteresis, also produces a variable damping, but with the damping becoming larger as the signal amplitude falls.

The nonlinear-gain element, Fig. 3(A), is certainly one possibility; others are discussed in subsequent sections. These latter fall into two groups: nonlinear rate feedback and nonlinear phase-lead compensation. Feeding back a signal proportional to some function of the system speed and introducing a cascade phase-lead network are widely practiced for improving a system's stability of response. Logically then, one or both of these tech-

Fig. 4. Schematic of system affected by coulomb friction



niques may be modified to suit the problem under consideration.

## Nonlinear Rate Feedback

### GENERAL FACTS

Some recent investigations on the effects of coulomb friction in feedback systems having hysteresis nonlinearities reveal that this friction has beneficial effects on stability.<sup>6,7</sup> This is because it produces, in effect, a nonlinear feedback of the rate signal. See Fig. 4. However, presence of coulomb friction has a generally bad effect on static accuracy and causes a zone of steady-state errors. Stabilization of feedback systems containing hysteresis by intentional introduction of coulomb friction may thus be undesirable, particularly in high-precision systems.

The good effects of coulomb friction on small-signal stability of a system can, however, be realized in practice if the secondary feedback loop seen in Fig. 4 is not caused by the friction but represents a suitable external feedback of the rate signal. This is explained in discussing the case of a simple second-order servo having backlash hysteresis.

### A NONLINEAR RATE FEEDBACK SYSTEM

Arrangements of the proposed nonlinear rate feedback system is shown in Fig. 5(A).

The secondary feedback loop, around the block representing the motor, results in the feedback of a nonlinear rate signal  $H$ , the form of which depends on the characteristics of the nonlinear block  $N$ . For the purpose of the present discussions, block  $N$  is supposed to have the characteristics shown in Fig. 5(B). This is realized in practice by placing a limiter ahead of a tachometer generator.

### ANALYSIS OF NONLINEAR RATE FEEDBACK SYSTEM

A frequency-response analysis of the system in Fig. 5(A) is possible by replacing nonlinearities with corresponding quasilinear describing functions. The equivalent system is shown in Fig. 6. The portion of Fig. 6 shown within the broken lines, representing the motor with the nonlinear rate feedback, can be replaced by a describing function  $G(j\omega, \hat{\theta}_m)$  given by

$$G(j\omega, \hat{\theta}_m) = \frac{K}{j\omega(j\omega T + 1)} \frac{1}{1 + \frac{kKN(\hat{\theta}_m)}{j\omega T + 1}} \quad (1)$$

$$= \frac{K}{j\omega[j\omega T + 1 + KkN(\hat{\theta}_m)]}$$

where  $K$  is the loop-gain,  $k$  is the feedback loop-gain,  $T$  is the motor time constant, and  $N(\hat{\theta}_m)$  is the describing function of the block  $N$  and is given by

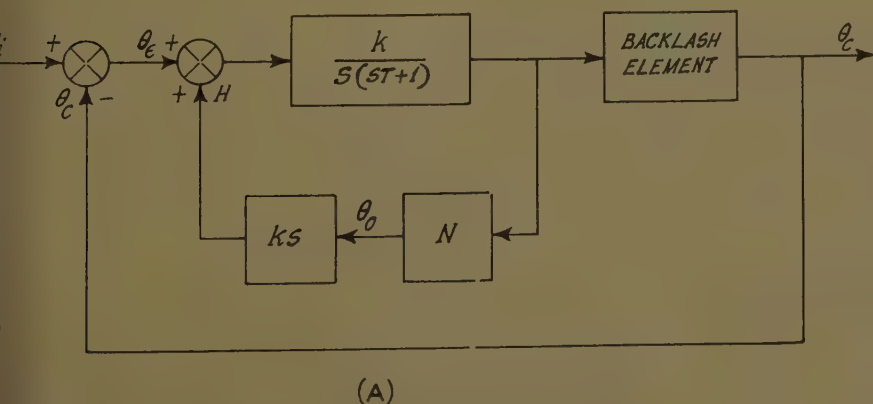


Fig. 5. A—Schematic of a nonlinear rate-feedback system. B—Characteristics of block  $N$

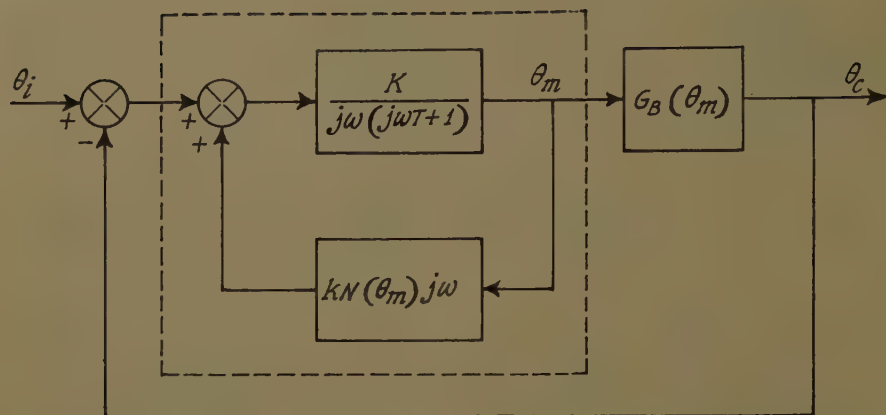


Fig. 6. Block diagram of system in Fig. 5 (A), using frequency-response functions

$$N(\hat{\theta}_m) = \frac{2}{\pi} \left[ \sin^{-1} \frac{s}{\hat{\theta}_m} + \frac{s}{\hat{\theta}_m} \sqrt{1 - \frac{s^2}{\hat{\theta}_m^2}} \right] \quad (2)$$

The simplified equivalent system can then be represented as in Fig. 7, and the condition for sustained oscillations is

$$G(j\omega, \hat{\theta}_m) = -[G_B(\hat{\theta}_m)]^{-1} \quad (3)$$

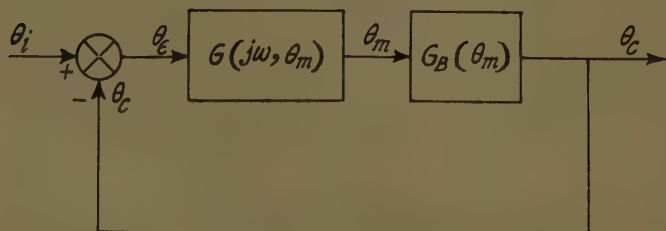
$G_B(\hat{\theta}_m)$  being the describing function of the backlash characteristics.

Validity of this condition in any system is tested most conveniently in a graphical manner, which also indicates suitable means for avoiding oscillations. The two functions in equation 3 are superposed, either in the complex (Nyquist or inverse Nyquist) plane or in the gain-phase shift plane. In Fig. 8, for example, are plots of these functions for assumed values of  $K=10$ ,  $k=0.1$ , and  $T=1$ ;  $n$  and  $m$  are normalized amplitude parameters, being respectively equal to  $\delta/\hat{\theta}_m$  and  $S/\hat{\theta}_m$ .

While the plot of  $-G_B(\hat{\theta}_m)^{-1}$  is an amplitude-dependent locus, those of  $G(j\omega, \hat{\theta}_m)$  are a family of frequency-dependent loci, each locus being drawn for a specific signal amplitude. Thus, validity of relation 3 requires that there be an intersection between the  $-G_B(\hat{\theta}_m)^{-1}$  locus and the particular  $G(j\omega, \hat{\theta}_m)$  locus, which bears the same amplitude mark as the  $-G_B(\hat{\theta}_m)^{-1}$  locus at the intersection point.

Fig. 7 (below). Equivalent representation of system in Fig. 6 for stability analysis

Fig. 8 (right). Nichols plots for system in Fig. 7 for  $K=10$  and  $T=1$



The forms of the  $G(j\omega, \hat{\theta}_m)$  loci in Fig. 8 reveal that, as a result of nonlinear rate feedback, the frequency-response transfer function of the motor unit has comparatively less gain and less lagging phase shifts at small-signal levels. This indicates a small-signal stabilizing effect of the nonlinear-rate feedback scheme of Fig. 5(A), and points out the possibility of avoiding sustained oscillations in the system by proper choice of the parameters of block  $N$ .

#### COMPENSATION FOR HYSTERESIS EFFECTS

Since the describing functions  $G(j\omega, \hat{\theta}_m)$  and  $G_B(\hat{\theta}_m)$  have been derived for  $\hat{\theta}_m(t)$ , which are assumed sinusoidal, the input signal  $\theta_i(t)$  is assumed to be a periodic function, making  $\hat{\theta}_m(t)$  vary sinusoidally. The relations then are as follows:

$$\theta_m/\theta_e = G(j\omega, \hat{\theta}_m)\omega$$

$$\theta_c/\theta_m = G_B(\hat{\theta}_m)$$

The total forward-loop transference is now

$$\theta_c/\theta_e = G(j\omega, \hat{\theta}_m) \cdot G_B(\hat{\theta}_m)$$

and the closed-loop transference is

$$\theta_c/\theta_i = G_B(\hat{\theta}_m) / [G(j\omega, \hat{\theta}_m)^{-1} + G_B(\hat{\theta}_m)] \quad (7)$$

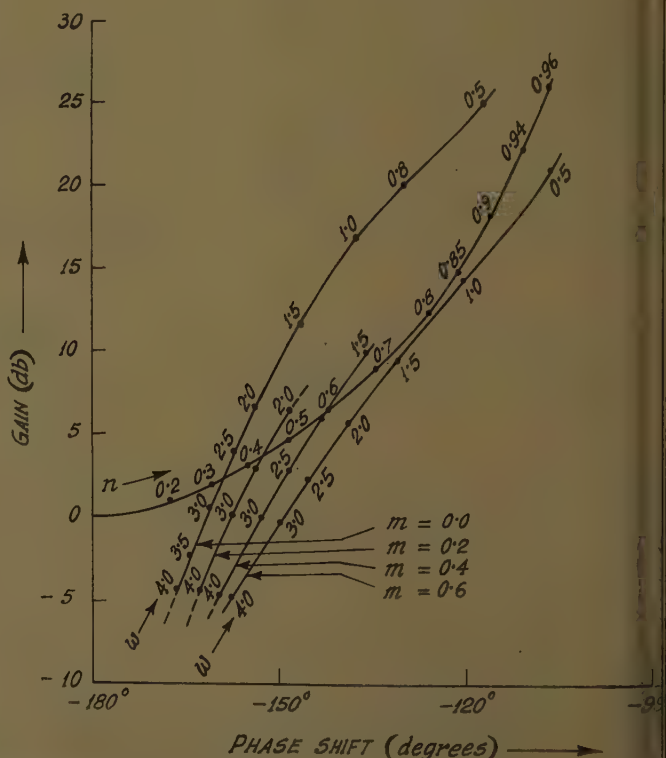
The frequency-response curves represented by equation 7 can be computed with the signal amplitude as a parameter by following a graphical procedure.<sup>10</sup>

Proper compensation for hysteresis effects requires that  $N(\hat{\theta}_m)$  be chosen in such manner that the amplitude dependence of these curves will be reduced to a minimum. This can be done by trial and error. But to design a compensated system by this procedure proves tedious and impractical. Adjustments are best made on the basis of computer runs and tests of a bread-boarded system.

#### Nonlinear Phase-Lead Compensation

##### GENERAL FACTS

The linear phase-lead network shown in Fig. 9(A) has the gain and phase functions shown in Fig. 9(B). Parameters of the network generally are so adjusted that the maximum phase lead occurs near the crossover frequency, resulting in improved phase margin with a consequent increase of relative stability. If such a network





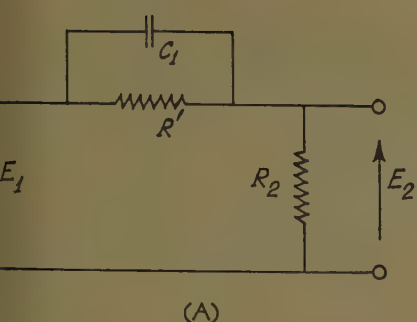


Fig. 9. A—RC phase-lead network. B—Bode plots of transfer gain and phase of the network

used to compensate for the phase-lag effects of hysteresis, the resultant system with greater phase margins at comparatively large signals implies that the system response should be somewhat sluggish for large-signal amplitudes. This apparent disadvantage is easily remedied by non-linearizing the network depicted in Fig. 9(A).

#### NONLINEAR PHASE-LEAD NETWORK

The network shown in Fig. 10(A) differs from the one in Fig. 9(A) in that two biased diodes have been added. The working of this circuit can be explained with the help of the equivalent circuit shown in Fig. 10(B) where  $R'$  is a nonlinear resistor, equivalent to the parallel combination of  $R_2$  and the diode conduction resistance. At small-signal levels, where  $E_s < E_1 - E_2 < E_s$ , both the diodes will be nonconducting and  $R' = R_2$ . The network then behaves in the usual manner. If, however, the signal amplitude is large enough to reverse the inequality, one

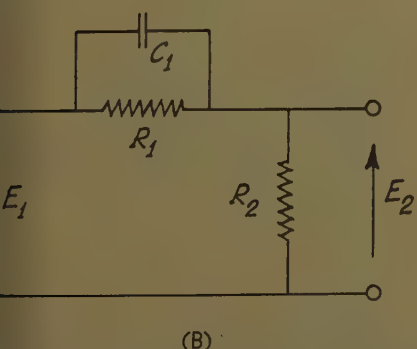
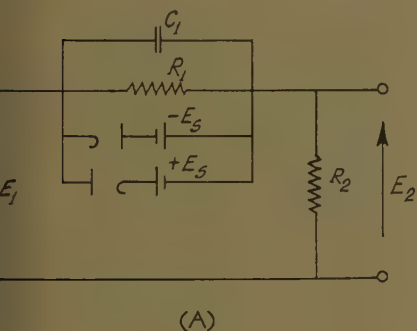
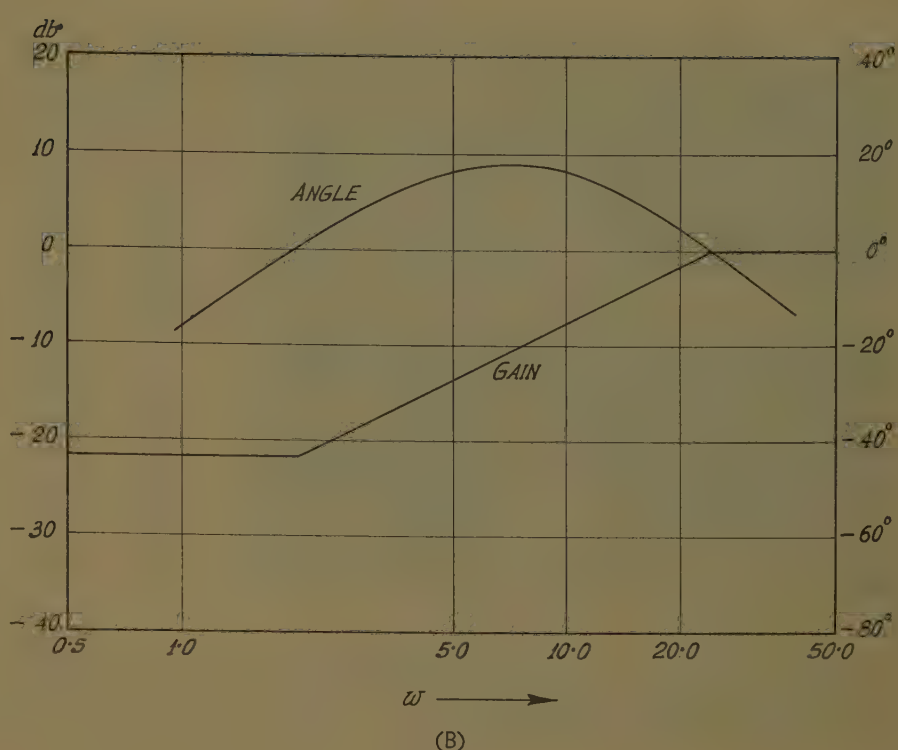


Fig. 10. A—Proposed nonlinear-lead network. B—Equivalent of the network



diode conducts, depending on signal polarity. Since diode conduction resistance is extremely small compared with  $R_2$ ,  $R'$  is also quite small and the circuit effectively behaves as a unity-transference network. The simple network of Fig. 10(A) thus is seen to introduce amplitude-dependent phase leads, which predominate at small-signal amplitudes only.

#### EFFECTS OF NONLINEAR LEAD NETWORK

The effects of incorporating the network of Fig. 10(A) in a nonlinear feedback system with hysteresis can be discussed on the basis of network frequency response. In Fig. 11, the hysteresis derives from backlash effects. The lead network is placed in tandem in the forward-loop after the error-sensing unit. This unit has a quasilinear frequency-response function  $G_c(j\omega, \hat{\theta}_e)$ , as indicated, which is a function of both frequency  $\omega$  and signal amplitude  $\hat{\theta}_e$ .

A derivation of this transfer function is difficult because conduction of the diodes depends on the signal difference  $E_1 - E_2$  ( $\theta_1 - \theta_0$ ) rather than on  $E_1$  ( $\hat{\theta}_e$ ) alone. (An easily workable, though approximate, describing function for the network is derived in the Appendix.

$$G_c(j\omega, \hat{\theta}_e) = [A_1^2 + B_1^2]^{1/2} / \tan^{-1}(B_1/A_1) \quad (8)$$

$$A_1 = \left[ \frac{q}{p} + \frac{2}{\pi} (1 - q/p) \cos^{-1} \frac{E_s}{\hat{E}_1} - \frac{2}{\pi} \frac{E_s}{\hat{E}_1} (1 + q/p) \sqrt{1 - \frac{E_1^2}{\hat{E}_1^2}} \right] \quad (9)$$

$$B_1 = \left[ \frac{\omega}{\pi p} \left( \pi - 2 \cos^{-1} \frac{E_s}{\hat{E}_1} + 2 \frac{E_s}{\hat{E}_1} \times \sqrt{1 - \frac{E_s^2}{\hat{E}_1^2}} - \frac{4q}{\pi p} \frac{E_s^2}{\hat{E}_1^2} \right) \right] \quad (10)$$

The best procedure under the circumstances is an experimental determination of transfer gain and phase as a function of signal frequency and amplitude.

Table I gives values of the instant network's transfer gain and phase for three different values of signal ratio  $E_1/\hat{E}_1$ , obtained on the basis of equation 8 and by direct measurement, the latter being within about 2% of absolute accuracy.<sup>11</sup> The order of approximation involved in equation 8 can now be estimated.

Function  $G_c(j\omega, \hat{\theta}_e)$  cannot be used as was  $G(j\omega, \hat{\theta}_m)$  for determining the limit cycle parameters, since the former is a function of  $\hat{\theta}_e$  while  $G_B(\hat{\theta}_m)$  is a function of  $\hat{\theta}_m$ . A trial-and-error procedure is permissible; for, under conditions of sus-

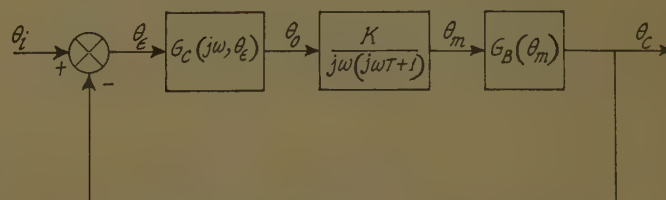


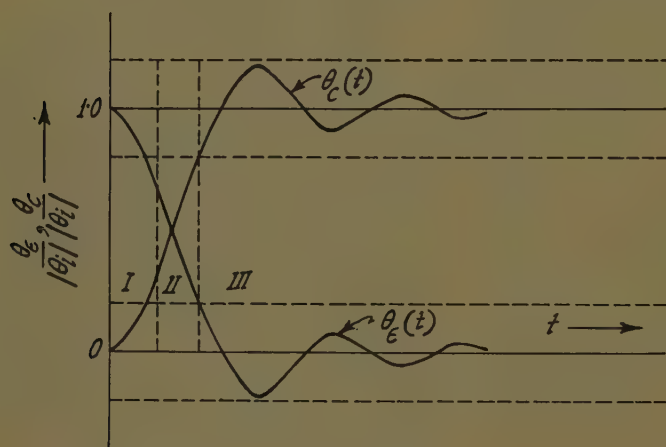
Fig. 11. Block diagram of the nonlinear lead-network compensated system

Table I. Transfer Gain and Angle of Nonlinear Lead Network

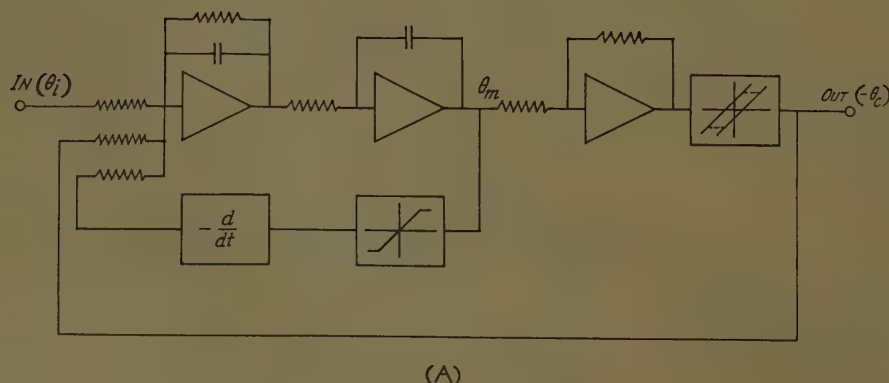
Three Regions In Network	Frequency, Radians Per Second	Gain		Angle, Degrees	
		Calculated	Measured	Calculated	Measured
I. $E_2/\hat{E}_1 = 1.0$ .....	1.86	0.114	0.115	38.5	38.5
	2.06	0.120	0.119	40.9	40.95
	2.62	0.136	0.139	46.3	46.4
	3.02	0.150	0.142	49.2	49.5
	3.59	0.169	0.173	52.4	52.4
	4.40	0.198	0.203	55.2	54.3
	5.70	0.245	0.250	57.4	56.4
II. $E_2/\hat{E}_1 = 0.8$ .....	7.85	0.321	0.329	57.6	56.7
	2.06	0.128	0.094	5.4	18.2
	2.62	0.131	0.096	13.9	22.3
	3.02	0.137	0.100	21.1	25.4
	3.59	0.145	0.104	28.3	28.8
	4.40	0.162	0.112	37.7	33.4
	5.70	0.194	0.125	49.1	39.4
III. $E_2/\hat{E}_1 = 0.4$ .....	7.85	0.264	0.149	61.0	48.0
	2.06	0.508	0.3856	3.0	3.05
	2.62	0.509	0.3870	4.2	3.8
	3.02	0.510	0.3872	5.3	4.5
	3.60	0.511	0.3875	6.5	5.3
	4.4	0.513	0.388	8.4	6.3
	5.7	0.518	0.390	11.3	8.4

Table II. Amplitude and Frequency of Oscillation in Nonlinear Rate Feedback System

K	Limiting Amplitude Compared with Backlash (s/s)	Values of $n = \delta / \hat{\theta}_m$		Values of Frequency, Radians Per Second	
		Simulator	Describing Function	Simulator	Describing Function
10	$\begin{Bmatrix} 0.0 \\ 0.5 \\ 1.0 \end{Bmatrix}$	$\begin{Bmatrix} 0.316 \\ 0.423 \\ 0.638 \end{Bmatrix}$	$\begin{Bmatrix} 0.315 \\ 0.440 \\ \text{No oscillations} \end{Bmatrix}$	$\begin{Bmatrix} 2.73 \\ 2.50 \\ 1.96 \end{Bmatrix}$	$\begin{Bmatrix} 2.80 \\ 2.30 \\ \text{No oscillations} \end{Bmatrix}$
4	$\begin{Bmatrix} 0.0 \\ 0.5 \end{Bmatrix}$	$\begin{Bmatrix} 0.601 \\ 0.732 \end{Bmatrix}$	$\begin{Bmatrix} 0.495 \\ \text{No oscillations} \end{Bmatrix}$	$\begin{Bmatrix} 1.3 \\ 1.02 \end{Bmatrix}$	$\begin{Bmatrix} 1.20 \\ \text{No oscillations} \end{Bmatrix}$



**Fig. 12. Typical error and response functions of system in Fig. 11, following step input disturbance**



tained oscillations, the assumption is that  $\hat{\theta}_1(t) = 0$ . Then  $\hat{\theta}_1(t) = -\hat{\theta}_c(t) = -G_B(\hat{\theta}_m(t))$ , and the condition of amplitude balance can be ascertained. However, such a medium is not necessary to gain understanding of the nonlinear network's effec-

Fig. 12 shows the form of type error and response functions of the system in Fig. 11 when  $\hat{\theta}_1(t)$  is a step displacement. The broken lines show the bounds of diode conduction. Network operation can then be explained for three different regions. Within the broken lines the network behaves as if diodes were absent and produces stabilizing phase lead effects of a conventional linear network. This is marked in Fig. 12 as zone III. For slightly larger signals, zone II in Fig. 12, the network introduces phase leads only slightly less than those in zone III, but its transfer gain is appreciably reduced (see data for  $E_s/E_1=0.88$  in Table I). This implies that the stabilizing action will here be greater for the network than for the conventional linear network. For even larger signals, zone I in Fig. 12, the network behaves almost as a unity-transference device and does not affect system operation.

Conclusions to be drawn from the foregoing observations are that, during initial stages of the system step response, little difference exists in the two networks when high-frequency transfer of unity is considered. The nonlinear network should give better response during later stages and should produce a greater stabilizing action. These conclusions are checked in the next section against the results of analog computer studies.

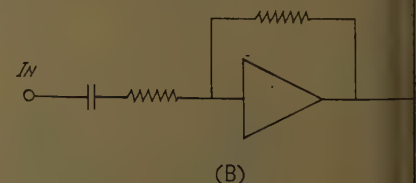
## Results of Simulator Studies

Experiments made with the help of analog simulators (real time) verified the preceding results, as described herewith.

## NONLINEAR RATE FEEDBACK SYSTEM

The simulator is arranged as shown in Fig. 13(A). The differentiator unit shown in Fig. 13(B) serves reasonably well in the frequency range of interest. Values of  $RC$  (resistance capacitance) elements are chosen to yield the transfer function

Fig. 13. A—Simulator arrangement for a linear rate feedback system. (B)—The derivative unit





$-10j\omega/(j\omega+100) \cong -0.1j\omega$  for  $j\omega < 100$ . First, the loop-gain is set high enough to yield sustained oscillations, the rate of change of the signal being set at zero. By fixing the saturation levels at different values, several simulated rate-feedback systems were made available for study. Table II contains amplitudes and frequencies of oscillations under different conditions. Values of these parameters, determined by using the describing functions, are also recorded, and they indicated the approximate accuracy to be expected of the describing functions.

The oscillograms in Fig. 14 display step response of the system for different settings of the limiter. The amplitude of the input step was set at  $\theta_1/\delta = 10$ . The limiter greatly increases the speed of response, while only slightly increasing the overshoot. The initial delay in traces Figs. 14(B) and 14(C) is expected, since full-rate signal is fed back during the initial part of the transient.

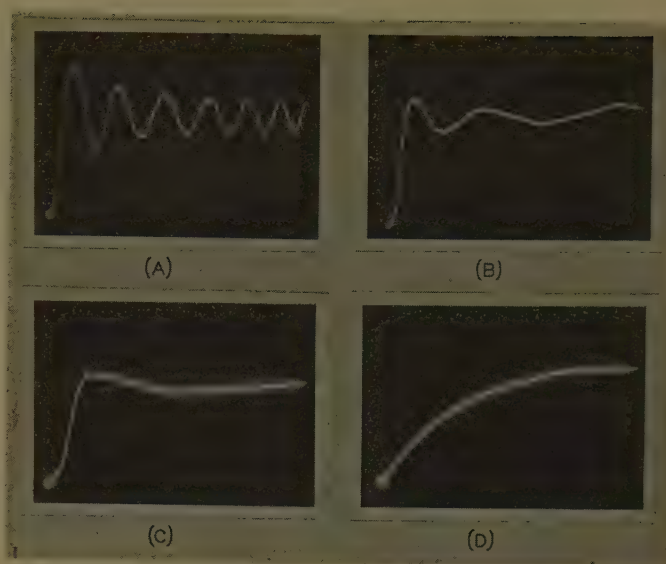
## NONLINEAR PHASE-LEAD SYSTEMS

The simulator setup is illustrated in Fig. 15, the compensating network being one shown in Fig. 10(A). The transfer function of the corresponding linear network—i.e., Fig. 10(A) with the diodes omitted—is  $(j\omega+2)/(j\omega+24)$ ; thus, the maximum phase shift occurs at about 1.1 rad (radians). Step responses of the system have been recorded, using two different values of the loop-gain:  $K=40$  and  $K=16$ . They give crossover frequencies of about 6.3 rad and 4.0 rad, respectively. Figs. 16(A) and (D) show the traces resulting from linear compensation in the two cases. Traces 16(B) and (E) result from nonlinearization of the compensation network (with the trace 16(A) superimposed by way of comparison) for  $E_s/\theta_1$  equal to 0.25 and 0.5, respectively, loop-gain being 40. When the value of the gain is 16, the traces are as shown in 16(F) and (F).

## CONCLUSIONS

To compensate for their undesirable effects on system stability, the hysteresis nonlinearities require nonlinear techniques. Usual rate-feedback and passive networks can be modified easily to help solve hysteresis problems. If a system incorporates a tachometer-feedback arrangement to achieve sufficient stability in the presumably linear system, then a linear-gain element (with larger gain for smaller signals) may be used in the feedback loop to guard against hysteresis effects. Also nonlinearization of the conventional RC phase-lead network, by the

**Fig. 14. Traces of step response obtained with setup of Fig. 13(A) where  $\theta_1/\delta = 10$ ,  $K = 10$ . A— $S/\delta = 0$ . B— $S/\delta = 1$ . C— $S/\delta = 2$ . D— $S/\delta = \infty$ .**



addition of biased diodes or otherwise, results in better response from the compensated system.

## Appendix

### Generalized Frequency-Response Function

When diodes are dismantled from the network seen in Fig. 10(A), the transfer function is

$$\frac{E_2(j\omega)}{E_1(j\omega)} = \frac{j\omega+q}{j\omega+p} = \frac{j\omega+2}{j\omega+24} \quad (11)$$

If the diodes are replaced, one of them will conduct for  $|E_1 - E_2|/E_s$ . If starting with  $E_1 = E_1 \cos \omega t$  so that at  $t=0$ , the lower diode with  $+E_s$  on the cathode will be conducting, and

$$E_2(t) = E_1(t) - E_s \quad (12)$$

the capacitor  $C_1$  being charged to  $+E_s$ . This continues until  $E_2(t)=0$  at  $t=t_1$ . After this instant, the effective circuit is the one shown in Fig. 17(A) which differs from the conventional linear-lead network by having the condenser charged to  $+E_s$ . The output voltage, after  $t=t_1$  is given by

$$E_2(t) = [E_1(t) - E_s] \left( \frac{s+q}{s+p} \right) \quad \left[ s = \frac{\alpha}{dt} \right] \quad (13)$$

This will continue until a time  $t=t_2$ , perhaps, when  $E_1(t) - E_2(t)$  become equal to  $-E_s$ , and the upper diode starts conduction. At  $t=t_2$

$$E_2(t) = E_1(t) + E_s \quad (14)$$

The switching instants  $t_1$  and  $t_2$  are illustrated in Fig. 17(B). The first instant is given by

$$\omega t_1 = \cos^{-1} \left( \frac{E_s}{E_1} \right) \quad (15)$$

To find the other instant, some approximations are made. Since  $p \gg q$ , the relationship, taken from equation 13, is

$$E_2(t) \approx \frac{E_1(t)}{p} (s+q) - \frac{E_s q}{p}$$

or from equation 14

$$E_2(t_2) - E_1(t_1) = \frac{E_1(t_2)}{p} (s+q-p) - \frac{q}{p} E_s = E_s$$

Thus

$$\frac{E_1(t_2)}{p} (s+q-p) = E_s (1+q/p)$$

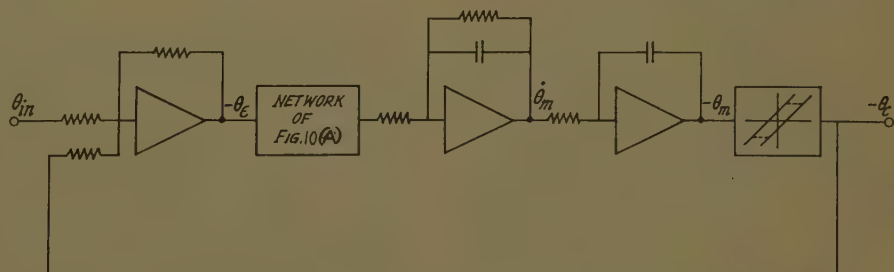
$$\text{or } \omega E_1 \sin \omega t_2 - (p-q) \cos \omega t_2 = E_s (p+q)$$

This gives for  $t_2$

$$\omega t_2 = \cos^{-1} \left[ -\frac{E_s}{E_1} \frac{p+q}{[(p-q)^2 + \omega^2]^{1/2}} \right] - \tan^{-1} \frac{\omega}{p-q}$$

$$\approx \pi - \cos^{-1} \frac{E_s}{E_1} \quad (16)$$

at low frequencies.



**Fig. 15. Simulator arrangement for the nonlinear lead-network compensated system**

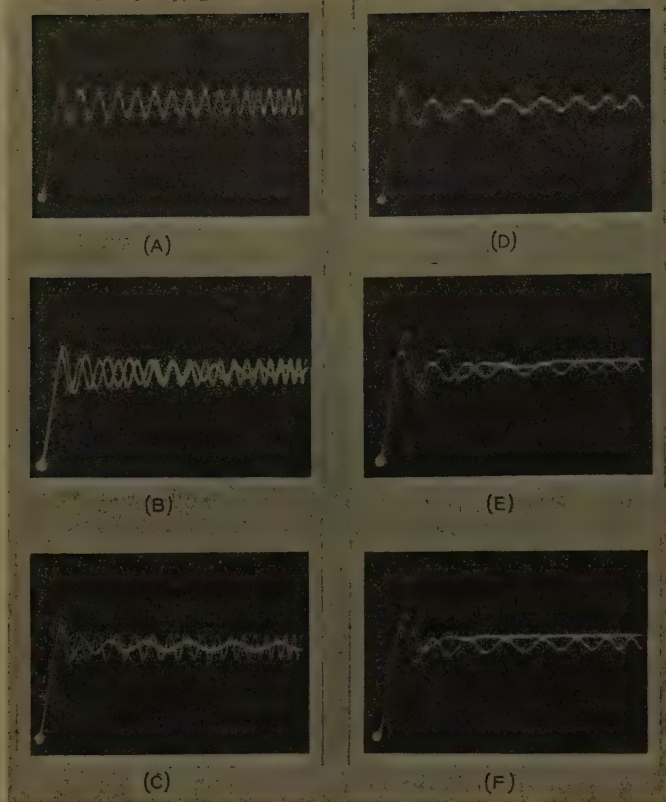


Fig. 16. Traces of step responses obtained with the setup of Fig. 15, where  $\theta_i/\delta = 10$ . A—K = 40,  $E_s/\theta_i > 1$ . B—K = 40,  $E_s/\theta_i = 0.25$ . C—K = 40,  $E_s/\theta_i = 0.5$ . D—K = 16,  $E_s/\theta_i > 1$ . E—K = 16,  $E_s/\theta_i = 0.25$ . F—K = 16,  $E_s/\theta_i = 0.5$

On carrying out the integrations

$$A_1 = \left[ \frac{q}{p} + \frac{2}{\pi} (1 - q/p) \cos^{-1} \frac{E_s}{E_1} - \frac{2}{\pi} (1 + q/p) \frac{E_s}{E_1} \sqrt{1 - \frac{E_s^2}{E_1^2}} \right] \quad (22)$$

$$B_1 = \left[ \frac{\omega}{\pi p} (\pi - 2 \cos^{-1} \frac{E_s}{E_1} + 2 \frac{E_s}{E_1} \sqrt{1 - \frac{E_s^2}{E_1^2}}) - \frac{4q}{\pi p} \frac{E_s^2}{E_1^2} \right] \quad (23)$$

## Nomenclature

- $A_1$  = real part of  $G_C(j\omega, \hat{\theta}_e)$
- $B_1$  = imaginary part of  $G_C(j\omega, \hat{\theta}_e)$
- $C_1$  = capacitor
- $D$  = width of magnetic hysteresis
- $E_1$  = lead-network input signal
- $E_2$  = lead-network output signal
- $E_s$  = diode bias
- $G_C(j\omega, \hat{\theta}_m)$  = describing function of motor with nonlinear rate feedback
- $G_C(j\omega, \hat{\theta}_e)$  = describing function of nonlinear lead network
- $G_B(\hat{\theta}_m)$  = describing function of backlash hysteresis
- $G_D(\hat{\theta}_e)$  = describing function of contactor hysteresis
- $G_H(\hat{\theta}_e)$  = describing functions of magnetic hysteresis
- $H$  = nonlinear rate signal
- $h$  = relay hysteresis
- $K$  = forward loop gain
- $k$  = feedback loop gain
- $m$  =  $S/\hat{\theta}_m$
- $N$  = nonlinear block
- $N(\hat{\theta}_m)$  = describing function of block  $N$
- $n$  =  $\delta/\hat{\theta}_m$
- $p$  = pole of  $RC$  network
- $q$  = zero of  $RC$  network
- $R_1, R_2, R'$  = resistances
- $s = d/dt$
- $S$  = limiting signal level
- $T$  = motor time-constant
- $\delta$  = backlash width
- $\Delta$  = dead zone
- $\theta_e$  = output signal
- $\theta_i$  = input signal
- $\theta_e$  = error signal
- $\theta_m$  = motor output signal
- $\theta_0$  = controller output signal
- $\hat{\theta}_m$  = amplitude of  $\theta_m(t)$  assumed sinusoidal
- $\hat{\theta}_e$  = amplitude of  $\theta_e(t)$  assumed sinusoidal
- $\theta_{es}$  = saturation level of magnetic hysteresis
- $\rho = D/\hat{\theta}_e$
- $\lambda = \theta_{es}/\hat{\theta}_e$
- $\omega$  = angular frequency

The network's generalized frequency-response function is now computed with the help of relations expressed in equations 12 through 16, giving

$$E_1 = \hat{E}_1 \cos \omega t \quad (17)$$

$$E_2 = \hat{E}_1 \cos \omega t - E_s, 0 \leq \omega t < \omega t_1 \quad (18A)$$

$$= \frac{\omega \hat{E}_1 \sin \omega t}{p} + \frac{q}{p} (\hat{E}_1 \cos \omega t - E_s), \quad \omega t_1 \leq \omega t \leq \omega t_2 \quad (18B)$$

$$= \hat{E}_1 \cos \omega t + E_s, \omega t_2 \leq \omega t \leq \pi \quad (18C)$$

Assuming

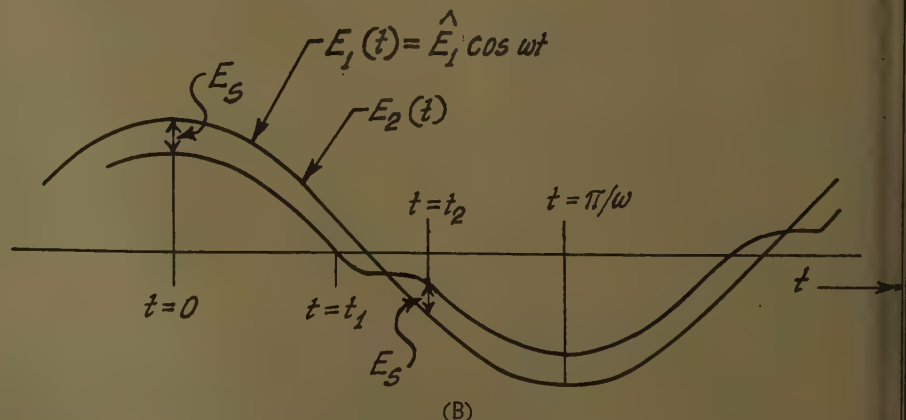
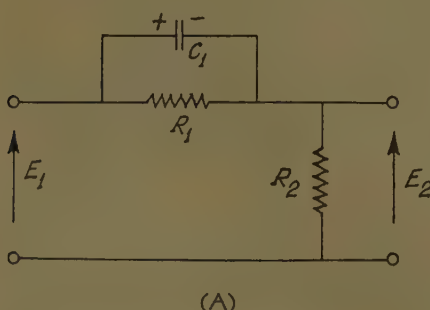
$$E_2/\hat{E}_1 = A_1 \cos \omega t_1 + B_1 \sin \omega t_1 + A_3 \cos 3\omega t + B_3 \sin 3\omega t + \dots \quad (19)$$

Then

$$A_1 = \frac{2}{\pi \hat{E}_1} \left[ \int_0^{\omega t_1} (\hat{E}_1 \cos \omega t - E_s) \cos \omega t \times d(\omega t) + \int_{\omega t_1}^{\omega t_2} \left\{ (\omega/p) \hat{E}_1 \sin \omega t + \frac{q}{p} (\hat{E}_1 \cos \omega t - E_s) \right\} \cos \omega t \alpha(\omega t) + \int_{\omega t_2}^{\pi} (\hat{E}_1 \cos \omega t + E_s) \cos \omega t d(\omega t) \right] \quad (20)$$

$$B_1 = \frac{2}{\pi \hat{E}_1} \left[ \int_0^{\omega t_1} (\hat{E}_1 \cos \omega t - E_s) \sin \omega t \alpha(\omega t) + \int_{\omega t_1}^{\omega t_2} \left\{ (\omega/p) \hat{E}_1 \sin \omega t + \frac{q}{p} (\hat{E}_1 \cos \omega t - E_s) \right\} \sin \omega t \alpha(\omega t) + \int_{\omega t_2}^{\pi} (\hat{E}_1 \cos \omega t + E_s) \sin \omega t \alpha(\omega t) \right] \quad (21)$$

Fig. 17. A—Equivalent of Fig. 10(A) for  $t_1 < t < t_2$ . B—Illustrative waveform of the nonlinear lead-network output defining the switching instants  $t_1$  and  $t_2$





## References

FREQUENCY RESPONSE METHOD FOR ANALYZING AND SYNTHESIZING CONTACTOR SERVOMECHANISMS, R. J. Kochenburger. *AIEE Transactions*, pt. I, 1950, pp. 270-84.

HARMONIC ANALYSIS OF FEEDBACK-CONTROL SYSTEMS CONTAINING NONLINEAR ELEMENTS, Johnson. *Ibid.*, pt. II (*Applications and Theory*), vol. 71, July 1952, pp. 169-81.

FEEDBACK IN CONTOURING CONTROL SYSTEMS, Ellert. *Ibid.*, vol. 74, 1955 (Jan. 1956 section), pp. 5-54.

NONLINEAR INTEGRAL COMPENSATION OF A VELOCITY-LAG SERVOMECHANISM WITH BACKLASH, Shen, H. A. Miller, N. B. Nichols. *Transactions*, American Society of Mechanical Engineers, New York, N. Y., vol. 79, 1957, pp. 585-92.

5. BACKLASH COMPENSATION IMPROVES SERVO PERFORMANCE, D. C. McDonald. *Instruments and Automation*, Pontiac, Ill., vol. 28, Nov. 1959, pp. 1728-31.

6. EFFECTS OF COULOMB FRICTION ON THE PERFORMANCE OF A SERVOMECHANISM HAVING BACKLASH, PART II, TRANSIENT RESPONSE CONSIDERATIONS, A. K. Mahalanabis. *Indian Journal of Physics*, Calcutta, India, vol. 35, no. 1, 1961, pp. 1-15.

7. STABILISATION OF CONTACTOR SERVOS BY USING COULOMB FRICTION, A. K. Mahalanabis. *Journal of Electronics and Control*, London, England, vol. 8, no. 4, 1960, p. 307.

8. NONLINEAR COMPENSATION NETWORKS FOR

FEEDBACK SYSTEMS, E. Mishkin, J. G. Truxal. *National Convention Record*, Institute of Radio Engineers, New York, N. Y., pt. 4, 1957, pp. 3-7.

9. A DESCRIBING FUNCTION STUDY OF HYSTERESIS EFFECTS IN A VELOCITY-LAG SERVOMECHANISM, A. K. Mahalanabis. *Proceedings*, First Defence Electronics Convention, Bangalore, India, Sept. 28-Oct. 1, 1959, pp. 99-103.

10. FREQUENCY RESPONSE OF NONLINEAR CLOSED-LOOP FEEDBACK CONTROL SYSTEMS, G. H. Fett, S. L. Mikhail. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 77, Nov. 1958, pp. 436-38.

11. A TRANSFER FUNCTION ANALYSER FOR LINEAR AND NONLINEAR COMPONENTS, A. K. Choudhury, M. S. Basu, A. K. Mahalanabis. *Electronic Engineering*, London, England, vol. 33, June 1961, pp. 382-85.

# Thermoelectricity Application Considerations

A. A. SORENSEN  
MEMBER AIEE

THESE DAYS, A DESIGNER is faced with myriad electric power requirements—from milliwatts to megawatts—dictated by new weapon and space systems concepts. At his disposal, fortunately, are numerous and varied potential sources of electricity. There is renewed interest in fuel cells and thermionic converters, for example. These, together with improved batteries and rotating generators driven by turbo prime movers, present many possibilities that warrant consideration for a power system. The particular selection, of course, is governed by an optimization process, covering required power level, operating time, and desired characteristics. The thermopile is a contender for many applications, now that solid-state physics advances have brought about more efficient uniform thermoelectric materials. It has both advantages and disadvantages. Compared with other energy-conversion devices, it is simple in construction, requiring no electrical fluids or extremely close spacings to maintain. Being static, the thermopile has the possibility of long life. It operates from chemical, solar, or nuclear heat sources, is able to use waste heat, and adapts to various operating conditions. On the other hand, its specific weight is high, voltage is low, and time constants are long. Even so, electromagnetic machines generate low voltage from the standpoint of individual turns. Similarly, in a thermopile series, connections are necessary to obtain required voltages. Power resistors and other solid-state devices efficiently convert direct current

to desired utilization levels and frequencies. Employing a storage device, long time constants can be handled more easily with d-c than with a-c sources.

## Engineering Design

Use of stable metals or alloys for thermocouples in measurement and control applications has grown into an important technology. Semiconductors also have been long recognized as powerful thermoelectric materials. Although only a few semiconductors are elemental in composition, there are numerous semiconducting compounds, many of which will give an electrical output when subjected to a temperature gradient. Some of these are particularly interesting because, having low thermal conductivity, they are economical in their use of heat. General requirements for practical use of these materials in generators are: high electrical output, low thermal conductivity, stability, and adaptability to production processes.

An important phase in any semiconductor materials program is control of the electron or hole concentration in the crystal. In elementary types, such as are found in diodes and transistors, this control is accomplished by adding the desired impurity atoms, which introduce donor or acceptor levels in the semiconductor's band structure. The situation is more complicated in binary compound semiconductors because the dominating carrier type is influenced by the deviation of the compound's composition from stoichiometry. Thus, the material's

properties are influenced by both composition and doping agent as well as the method of preparation.

One such thermoelectric material is lead telluride. Couples of the highest efficiency are obtained if *n*-type material is prepared with excess lead over the stoichiometric ratio and *p*-type material with excess tellurium. However, since material with excess lead is mechanically superior to excess tellurium material, it lends itself better to mechanical processes necessary in the construction of a generator.

The preparation method also affects mechanical characteristics. For example, pressed and sintered lead telluride is machined more readily than cast material. Pressing and sintering also contribute to preparing graded thermoelectric arms; that is, to greater doping of the hot end of the arm than the cold end. This is used because the temperature depends on the material's properties and results in some increase in efficiency. Characteristics of graded *n*- and *p*-type arms, produced by pressing and sintering, are shown in Figs. 1 and 2.

More than optimum *n*-type and *p*-type thermoelectric arms are required in engineering applications, however. Also required are a heat source, electrical insulation for hot and cold sides of the couples, thermal insulation, electrical connectors for hot and cold sides of the arms, and a heat sink. As far as the heat source and sink are concerned, the problem in engineering design is to match available heat fluxes with those required for thermoelectric elements. Probably the most troublesome item is joining electrical connections to the thermoelectric arms. Soldering is generally adequate for

Paper 60-1066, recommended by the AIEE Aerospace Transportation Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE Pacific General Meeting, San Diego, Calif., August 8-12, 1960. Manuscript submitted March 23, 1959; made available for printing May 16, 1961.

A. A. SORENSEN is with The Martin Company, Baltimore, Md.

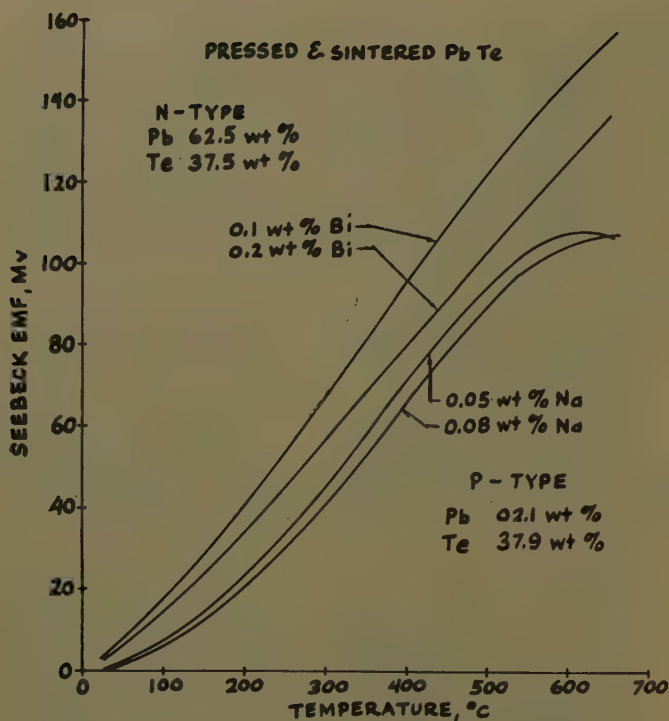


Fig. 1. Seebeck electromotive force versus temperature

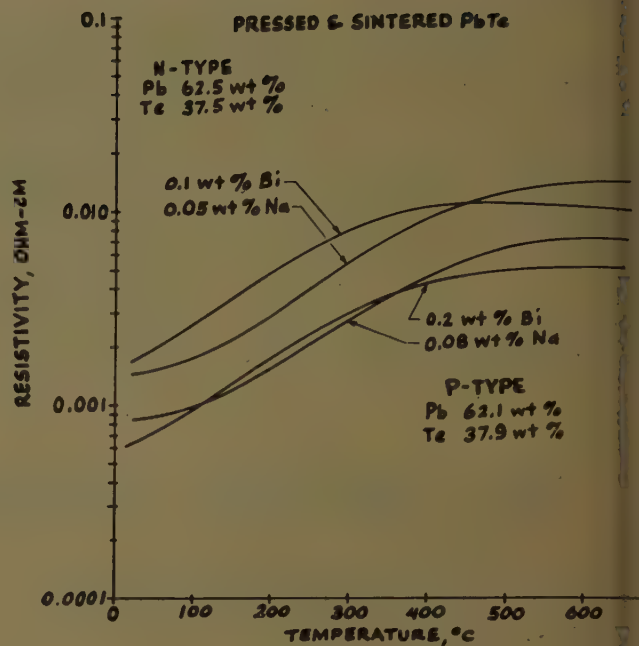


Fig. 2. Resistivity versus temperature

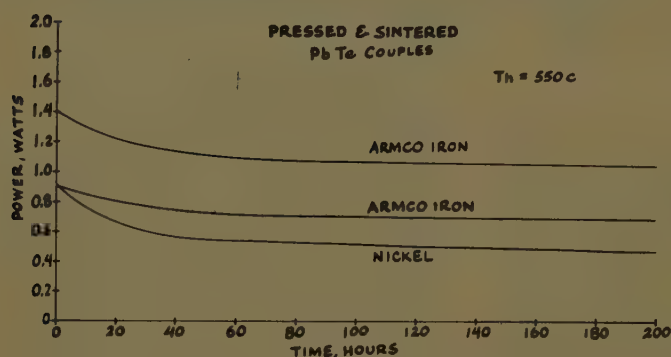


Fig. 3. Power output versus time

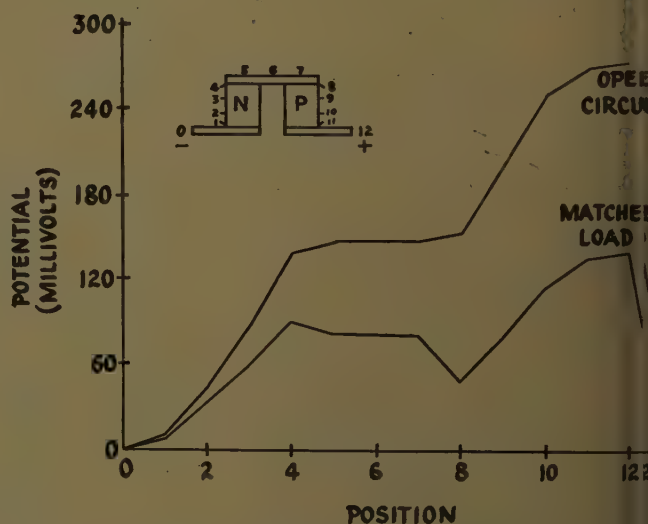


Fig. 4. Couple potential profile

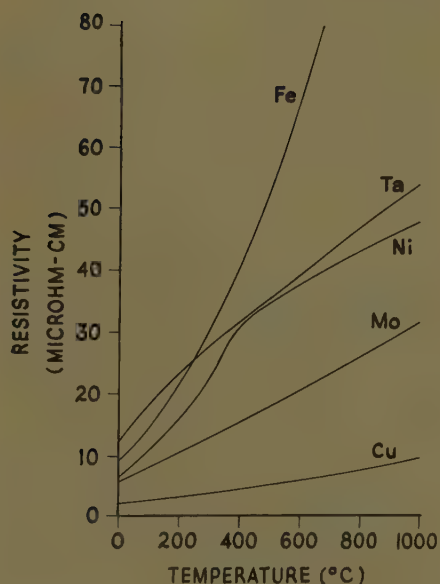


Fig. 5. Resistivity of some metals

Fig. 6 (right). Figure of merit versus temperature for typical thermoelectric materials

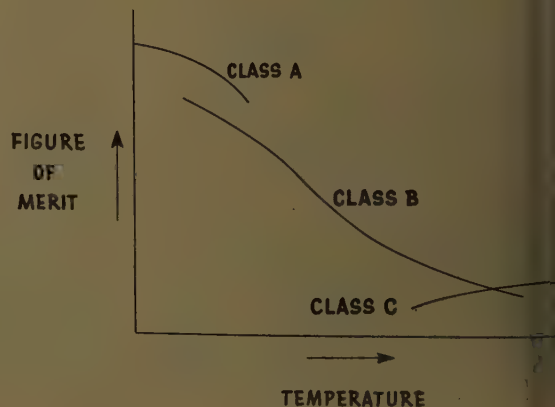




Table I. Space Applications

Applications	Ratings	Limitations
Electric power source for rocket stages, using cryogenic cooling and rocket at	1 to 20 kw (kilowatts) (15 to 25% efficiency)	Weight, thermal lag
Electric power source for satellites, using nuclear or solar heat sources	1 watt to 5 kw (5 to 8% efficiency)	Pumping losses limit efficiency intermediate ratings. Heat storage needed for solar source
Control power devices	Milliwatts to watts	

Table II. Terrestrial Applications

Applications	Ratings	Limitations
Electric power source for automotive equipment	5 to 10 kw (10% efficiency)	Best suited for low speed, intermittent use. Battery needed for peak loads
Combustion-powered, portable electric power source for quiet operation	100 watts to 5 kw (5 to 10% efficiency)	Pumping losses an appreciable part of output
Combustion-powered battery chargers for use in remote areas	25 watts to 1 kw (2 to 7% efficiency)	Orientation and concentration required for higher efficiencies

cold sides. For the hot sides, however, materials must be used which do not affect the performance of either *n*- or *p*-type arms. Fig. 3 shows the output of three couples made with lead telluride, one with iron electrodes, and one with nickel. The drop-off in output with time is more rapid for the nickel electrode be-

cause the hot end of the *p*-type material changes to *n*-type.

A helpful tool for gauging over-all performance of a thermocouple is the couple-potential profile shown in Fig. 4. Fine wire voltage pickups are attached at points on the thermocouple as indicated by the numbers shown. Voltages across

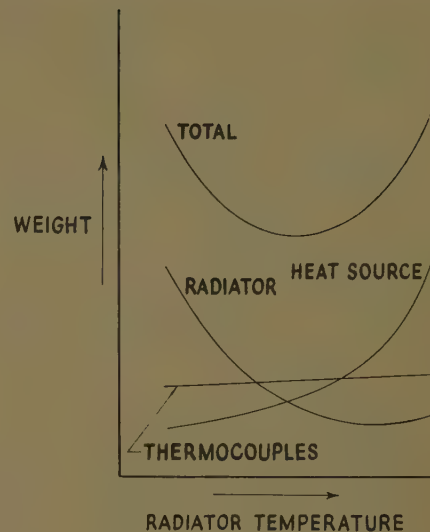


Fig. 7. Graphical solution for minimum weight for heat source, thermocouples, and radiator

cold-end connections are indicated by 0 to 1 and 11 to 12; across the hot-end connections by 4 to 5 and 7 to 8. The open-circuit profile of the couple shows voltage distribution between *n*- and *p*-type arms. This also reveals deterioration in performance during extended operation and points out the cause of the difficulty, such as reverse doping of an arm by the electrode. In this case, a reverse voltage is seen at the hot end of the arm involved. The matched-load profile permits determination of contact resistances under operating conditions and shows output distribution between the two arms. The profile also allows voltage drop in the body of the electrode connecting the two arms to be checked.

The hot connecting electrode can be-

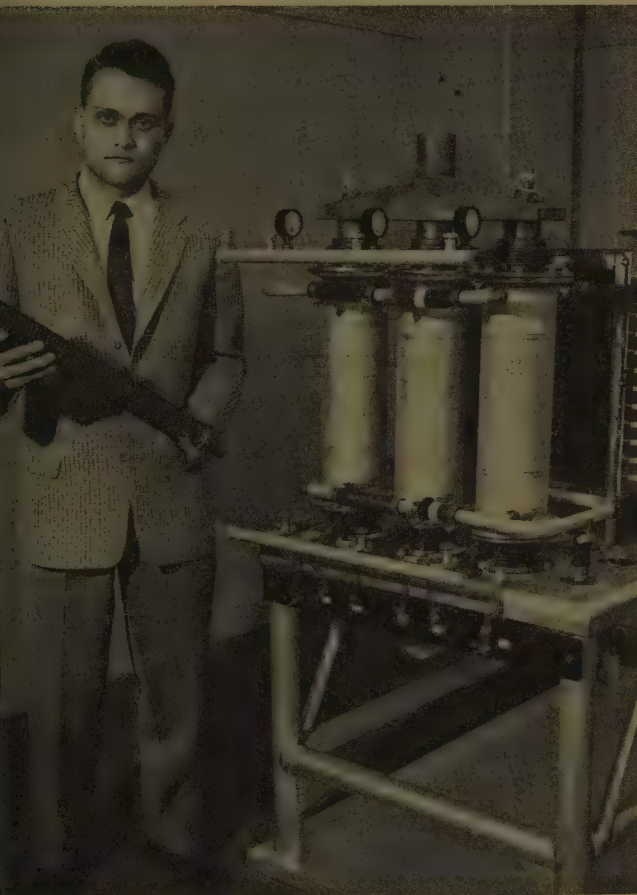


Fig. 8. Thermoelectric generator

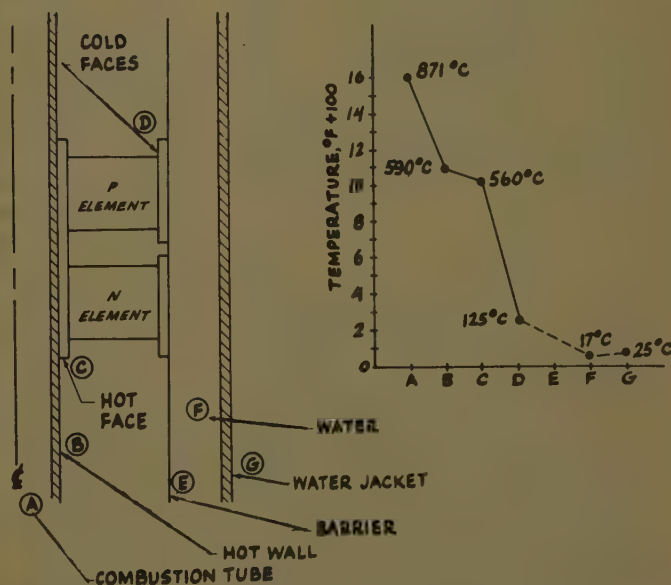


Fig. 9. Temperature profile

Table III. Marine Applications

Applications	Ratings	Limitations
Main propulsion	500 kw and up (15 to 20% efficiency)	Best used with battery for system voltage control
Auxiliary power	250 watts to 100 kw (5 to 20% efficiency)	
Signal devices	5 to 50 watts (10 to 15% efficiency)	

Table IV. Net Thermocouple Efficiency

Applications	Temperature		Efficiency		
	$T_h$	$T_c$	$Eff_o$	$Eff_{tc}$	$Eff_N$
Space	2,000 degrees Fahrenheit (1,093 degrees centigrade)	900 degrees Fahrenheit (482 degrees centigrade)	45%	17%	7.7%
Terrestrial	1,600 degrees Fahrenheit (871 degrees centigrade)	150 degrees Fahrenheit (66 degrees centigrade)	70%	26%	18%
Marine	1,600 degrees Fahrenheit (871 degrees centigrade)	45 degrees Fahrenheit (7 degrees centigrade)	75%	28%	21%

come an appreciable weight item, particularly when contact resistance is lowered and the element lengths decreased. Fig. 5 shows resistance versus temperature for some of the metals of interest for contacts. At higher temperatures, composite connectors are indicated to obtain minimum connector volume.

Efficiency is dependent upon the figure of merit (usually designated by  $Z$ )—a composite of the Seebeck coefficient, thermal conductivity, and resistivity. Thus far, the highest figures of merit have been attained by materials suitable for use in room-temperature range. The general trend as temperature increases is shown in Fig. 6. Representative of class *A* are various alloys of  $Bi_2Te_3$ . Class *B* includes  $PbTe$  and  $LaSb$ , while class *C* covers oxides and silicides. Conditions under which design-work operation is to take place must be established—the heat sink temperature, range of temperature from the heat source, required life, and characteristics required by the load. From the standpoint of effect upon the thermopile, designs fall into three general categories: space, terrestrial, and marine.

#### SPACE

From the designer's viewpoint, space is the worst category because radiation is the only heat-dissipating agent available. For small power applications, conduction may be used to transport heat to the radiator. If the power is appreciable, though, weight becomes excessive for this method, and a circulating fluid system must then be used. When a design is optimized for weight by compromising between the radiator's weight and that of the thermopile and heat source, then the radiator operates at a comparatively high

temperature. Optimization may be carried out by analytical solution of equations, by a computer, or by graphical solution based on calculations of several cases, as shown in Fig. 7.

Since the temperature gradient in the elements is assumed to be fixed, the weight of the thermoelectric elements will not change greatly with the radiator temperature. Thus, for a fixed maximum, the greater number of elements required at a higher radiator temperature is compensated by their shorter length. This assumption requires that contact resistance must not be the governing factor in element length. Table I shows some possible space applications.

#### TERRESTRIAL

Although nuclear, solar, and combustion heat sources are available for terrestrial use, the latter likely will provide for the majority of cases. Catalytic combustion processes are especially favorable for thermoelectric generation, since an extended high-temperature area can be provided and heat transfer is largely by radiation, thus reducing stack losses. Fig. 8 shows a generator employing this combustion principle; the catalytic combustion tube held by the young man is

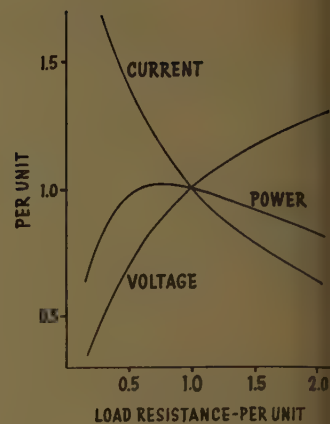


Fig. 10. Output parameters versus load

similar to ones inside each of the conversion tubes. At the time the photograph was taken, the output of the generator with pressed and sintered telluride couples was 130 watts. Each of the stacks weighed  $46\frac{3}{4}$  pounds. Fig. 9 shows a profile of temperature distribution from heat source to sink.

For terrestrial uses, as contrasted with space, both radiation and convection effects will dispose of unusable heat. Thus, radiator weight can be reduced. However, pumping power is needed to move air through a radiator, and an optimization must be made between over-all efficiency and weight. Military specifications for ground applications require the maximum cooling-air temperature be 125 degrees Fahrenheit, and lower, depending upon the use, for airborne classifications. Table II shows possible terrestrial examples.

#### MARINE

Marine applications, some of which are outlined in Table III, are the most favorable from a heat-disposal viewpoint because water is available in quantity at a temperature no higher than 80 degrees Fahrenheit. Thus, the machinery for heat disposal is simple and the pumping power requirements are a small percentage of the total power generated.

The design of a complete power-conversion system is governed not only

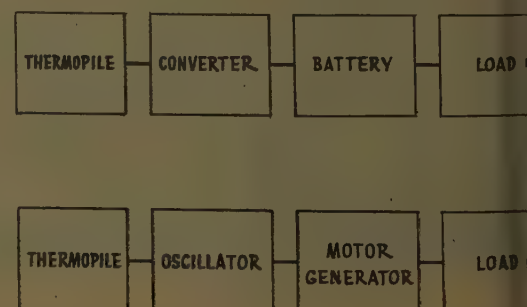


Fig. 11. Arrangements for supplying cyclic loads



capabilities of thermoelectric materials but is also strongly influenced by temperature limitations of the heat source, electrodes, seals, support materials, and electrical and thermal insulation. Although the marine application is the most favorable from the efficiency standpoint, it does not infer that it alone is practical.

In each area, thermoelectric power systems are subject to the same engineering-material limitations as competing systems, except for one advantage: they do not withstand the high stresses encountered with rotating systems. The relative advantage for operation at higher temperatures in space applications will demand more time and effort in materials development and engineering effort. A comparison of near-term thermocouple efficiencies for the three areas is shown in Figure IV. Over-all equipment efficiencies will be lower because of the factors previously discussed.

## Methods of Use

With currently possible efficiencies of about 10%, the value of  $m$  (ratio of external to internal resistance) is approximately 1.3. This means that, at the maximum efficiency point, the ratio of open-circuit to closed-circuit voltage with rated load is 1.77, or a regulation of 43.5%, which is equivalent to that of an unregulated a-c generator. With improved materials,  $m$  will be larger and the regulation lower in value. Voltage, current, and power are plotted against per-unit resistive load in Fig. 10. Note that power output is nearly constant for a wide range of load resistance. Applications requiring a constant power input, such as some motor drives, could take advantage of this characteristic. In other applications, needing direct current at various voltages, transistor converters with internal voltage regulators

may be used. Similarly, for a-c needs, inverters are chosen.

Thermopiles have the important function of supplying continuous power to a storage device, such as a battery or rotating flywheel, from which energy is withdrawn as needed. In this application, the thermopile can be operated at its most efficient load. Intermediate equipment performs the task of regulating the voltage and converting it to the required form. Two arrangements are shown in Fig. 11.

Each source of electric energy—be it an electromagnetic generator, primary battery, or thermopile—has its distinctive characteristics and limitations, which engineers recognize and take into account by fitting designs to the various idiosyncrasies. With the thermopile becoming generally available, engineers will become adept in designing around its limitations and exploiting its advantages.

# Some Recent Advances in the Analysis and Synthesis of Nonlinear Systems

ALFRED A. WOLF  
MEMBER AIEE

DURING the past two decades a great deal of material has appeared in the literature concerning the analysis and synthesis of nonlinear systems. In reference 1, a great many of these techniques are reviewed and it is quite evident that the methods are highly specialized, pertaining to particular problems. Some linearization techniques and others are approximations, sometimes quite restricted. Each has its particular value and usefulness. Despite the number of methods, there seems to be no underlying unifying theory connecting them.

Very recent publications appearing in the literature<sup>2-8</sup> have shown the development of a mathematical theory for the analysis of a class of nonlinear systems. Beginning with the author's dissertation,<sup>2</sup> which gives a good account of this theory, series of articles have appeared expounding the deterministic process of analysis of nonlinear systems and extending the theory to systems subjected to stochastic processes with the aid of the Statistical Transform Theorem.<sup>9,10</sup> Morse,<sup>11</sup> working in collaboration with Ku,<sup>12,13</sup> and Lee,<sup>14</sup> has developed a linear theory from a point of view

quite different from that developed by the author in collaboration with Ku and Dietz. The Wiener-Lee-Bose theory emphasizes a statistical approach to the synthesis of nonlinear systems taking advantage of orthogonal networks such as the Laguerre and Hermite systems.<sup>15</sup>

It is the purpose of this paper to develop new theoretical relations between the statistics of the output of a system, when subjected to a white-noise probe, and the configuration of the topology of the system's actual structure. As a direct implication of this last point, a relation exists between theoretical facets of analysis and the corresponding aspects of synthesis. In addition, as a point of departure, the partition theory will be discussed and extended.<sup>2,3</sup> It will be shown that although the solution of a certain class of nonlinear systems is unique, the form of solution is not unique and is a function of the point of partition. Examples will be given to illustrate this. Finally, consideration will be given to a synthesis procedure for specified waveforms.

Other possible approaches to the solution of nonlinear systems will not be

considered since they are adequately treated in reference 1.

## Partition Theory Analysis

The systematic study of physical nonlinear systems is based on a study of a certain class of nonlinear differential equations. Consider a typical equation of this class:

$$Z(D)x(t) + F\{x(t), \dot{x}(t), \dots\} = g(t) \quad (1)$$

where the linear part  $Z(D)x(t)$ , the nonlinear part  $F\{\dots\}$ , and the driving functions  $g(t)$ , are suitably restricted by the broad class of conditions given in Appendix I. A result of imposing these conditions is that the solution is unique and analytic. Keeping these two important

Paper 61-713, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE-AIEChE-ASME-IRE-ISA Joint Automatic Control Conference, Boulder, Colo., June 28-30, 1961. Manuscript submitted October 26, 1961; made available for printing May 8, 1961.

ALFRED A. WOLF is with Emertron, Inc. (a subsidiary of Emerson Radio and Phonograph Corporation), Silver Spring, Md.

The author would like to thank Dr. Y. H. Ku of the University of Pennsylvania whose helpful discussions and comments were much appreciated, Dr. L. F. Kazda of the University of Michigan who suggested and encouraged this study, Dr. T. L. Higgins of the University of Wisconsin for his encouraging remarks on this work, Dr. J. Randolph, Dr. D. Shen, Dr. L. Kanal, and Mr. J. Dietz for many hours of stimulating discussion. The author is also grateful to Mr. K. Lord, Dr. R. Weller, and Mr. D. Keim of the Stromberg-Carlson Company for their encouragement, and finally Mr. Parkhill and Mr. Rapuano for their comments. This paper is a modified version of CP60-110 presented at the AIEE Winter General Meeting, New York, N. Y., February 4, 1960.

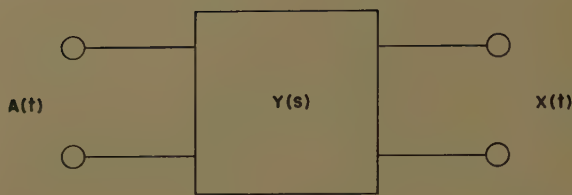


Fig. 1. A linear system with response  $x(t)$

results in mind, consider another problem. Referring to Fig. 1, let  $Y(s)$  define a transfer function of a linear system related to the linear differential operator  $Z(D)$  according to the relation

$$Y(s) = \frac{1}{Z(s)} \quad (2)$$

This is equivalent to the relation<sup>16</sup>

$$y(t) = \frac{1}{2\pi j} \int_{Br_1} \frac{e^{st}}{Z(s)} ds \quad (3)$$

by noting it is the inversion of equation 2 in the complex  $s$ -plane along a Bromwich contour ( $Br_1$ ) enclosing the poles of  $Z^{-1}(s)$ . An interesting question arises: Does a function  $A(t)$  exist such that when it is applied to the input terminals of  $Y(s)$ , shown in Fig. 1, a response  $x(t)$  results identical with the response of a nonlinear system driven by  $g(t)$  described in equation 1? The answer, from strictly physical considerations, is evidently in the affirmative. That this is true is clear at once if the nonlinear system is restricted to physical systems. Then  $x(t)$  is physically realizable.<sup>17</sup> Hence,  $A(t)$  is physically realizable and therefore it does exist. If function  $A(t)$  were known, it would be an easy task to calculate  $x(t)$ , so that  $x(t)$  would be determined by the familiar convolution integral:

$$x(t) = \int_0^t A(\tau) y(t-\tau) d\tau \quad (4)$$

Utilizing the conditions of Appendix I, the author showed in references 2 and 3 a rigorous mathematical basis for the existence and uniqueness of  $A(t)$ , the auxiliary forcing function. Moreover,  $A(t)$  is analytic under these conditions.

A hint as to how to calculate  $A(t)$  is obtained from physics;<sup>18</sup> consider Fig. 2. This system will be called the canonical system since it can be described by equation 1 which is the general form for a wide variety of physical systems satisfying conditions given in Appendix I. Suppose now that this system is partitioned as shown in Fig. 3 by making cuts at points  $a$  and  $b$ , shown in Fig. 2. By properly maintaining the system dynamics, the response of the linear part behaves exactly as the response of the interconnected system.

The equations describing the partition in Fig. 3 are

$$Z(D)x(t) = A(t) \quad (5)$$

and

$$A(t) = g(t) - F\{x(t)\} \quad (6)$$

Equation 6 is known as the auxiliary equation<sup>2,3</sup> while equation 5 is called the partition equation. It is evident that these equations could have been derived directly by mathematical manipulation of equation 4 with equation 1. What do these equations mean? Obviously, eliminating  $A(t)$  from equations 5 and 6 does not achieve any useful result except to give back equation 1. To use equations 5 and 6 effectively, advantage must be taken of the a priori properties of  $A(t)$ , which were deduced by means of the conditions required of the linear and nonlinear part of equation 1 and its forcing function. From these conditions,  $x(t)$  is deduced to be analytic. Since  $g(t)$  is analytic and  $F$  is continuous ( $A(t)$  is analytic since it is the difference between two analytic functions), the function  $F$  is taken to have the form

$$F(x) = \sum_{k=2}^N a_k x^k \quad (7)$$

This is not a necessary condition. If  $F$  does not have the form of equation 7, it must then be analytic. This condition may, in certain circumstances, be relaxed with due caution.

Thus there is justification for expanding  $A(t)$  into a power series. Under these

conditions the solution of equation 1 is given by the moment theorem.

## MOMENT THEOREM

Given an ordinary nonlinear differential equation (such as equation 1 suitable restricted according to the conditions given in Appendix I), the solution  $x(t)$  is a linear combination of all the moments of the folded impulse response of the linear part  $Z(D)$  where the coefficients are given as a recurrence equation.<sup>2-5</sup>

Formally, the moment theorem can be written as follows. Let  $Q_n(t)$  denote the  $n$ th moment of the folded impulse response of the linear part of equation 1 on the half-open interval  $(0, t)$ , i.e.,

$$Q_n(t) = \int_0^t \tau^n y(t-\tau) d\tau$$

Then the solution of equation 1 is:

$$x(t) = \sum_{n=0}^{\infty} c_n Q_n(t)$$

The coefficients  $C_n$  are determined recursively. When  $F(x)$  is a polynomial that given by equation 7 and the system is initially at rest, i.e., all initial conditions are zero, the recurrence equation for equation 1 is given by equation 10:

$$C_n = b_n - \sum_{k=2}^M C_{n-\alpha}^{(k)} a_k; \alpha \geq 0 \text{ and integral valued}$$

where the  $b_n$ 's are the Taylor coefficients of the expansion of the forcing function  $g(t)$ . Equation 10 is obtainable directly by the Taylor-Cauchy transform.<sup>8</sup>

## RULES FOR PARTITION

1. The dynamical equation representing the system is partitioned so that the linear terms involving the highest-order

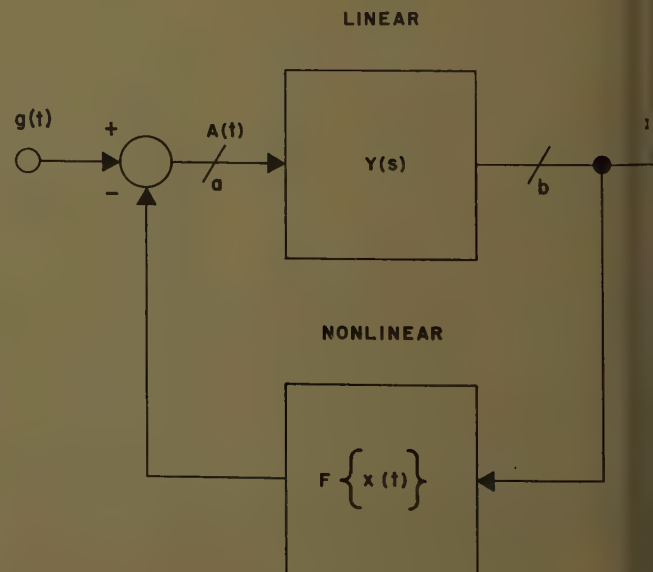


Fig. 2. A canonical nonlinear feedback system



itive are contained in one member equation.

A modified forcing function  $A(t)$  is obtained which consists of the actual g function, nonlinear and possibly linear terms.

A proper choice is made of the ex- on of  $A(t)$  according to the particular rties of the linear, nonlinear, and g functions. An auxiliary equa- thus formed.

The properly partitioned differen- quation is then solved.

Certain coefficients which arise the auxiliary equation are then cal- d.

These coefficients, which appear in partitioned equations, are then ated to obtain the exact solution.

## Taylor-Cauchy and Generalized Transforms

en the partition is made such that ighest-order derivative remains by in left member, the moment func-  $Q_n(t)$  is a power function. Hence lution,  $x(t)$ , is a power series. The ence relations for this partition are able by means of the Taylor- ay transform given by equations 11 2:

$$\frac{1}{2\pi j} \int_C \frac{W^{(k)}(\lambda)}{\lambda^{n+1}} d\lambda \quad (11)$$

$$x(\lambda) = \sum_{n=0}^{\infty} w_{n,k} \lambda^n \quad (12)$$

transform pairs  $w_{n,k}$  and  $W^{(k)}(\lambda)$  may lated operationally by the Taylor- ay operator  $\mathcal{J}_c$  according to the ons

$$\mathcal{J}_c \{ W^{(k)}(\lambda) \} \quad (13)$$

$$x(\lambda) = \mathcal{J}_c^{-1} \{ w_{n,k} \} \quad (14)$$

direct transform of  $W^{(k)}(\lambda)$   $\lambda) = k$ th derivative of a complex time function  $\lambda$

l discrete variable in the transform half line corresponding to the complex time variable in the complex time plane

ircle in the  $\lambda$  plane defining the domain of definition of the complex time function  $W^{(k)}(\lambda)$  and enclosing the singularities of the integrand of equa- tion 11

complex time variable corresponding to the real variable  $t$

ations 11 and 12 are used with a of Taylor-Cauchy transforms to nonlinear differential equations in ner analogous to the use of Laplace orms with linear constant parameter

differential equations. Details of this transform are given in reference 8 with adequate illustrations.

Consider now the conditions for the generalization of this transform to other partitions. Starting with equation 9 and again generalizing the real variable  $t$  to its analytic continuation in the complex  $\lambda$  plane, equation 15 is obtained:

$$x(\lambda) = \sum_{n=0}^{\infty} C_n Q_n(\lambda) \quad (15)$$

where

$x(\lambda)$  = a complex function derived from the real function  $x(t)$  by analytic continuation just as  $W(\lambda)$  was obtained from  $x(t)$  in the case of the Taylor-Cauchy transform

$C_n$  = corresponding transform mate of  $x(\lambda)$

Multiplying equation 15 by  $Q_{-m-1}(\lambda)$  and integrating around a contour  $C$  which encloses the appropriate singularities of the integrand gives interchanging orders of summation and integration.

$$\frac{1}{2\pi j} \int_C x(\lambda) Q_{-m-1}(\lambda) d\lambda = \sum_{n=0}^{\infty} C_n \frac{1}{2\pi j} \times \int_C Q_n(\lambda) Q_{-m-1}(\lambda) d\lambda \quad (16)$$

where the interchange is justified since the sum is uniformly convergent, under the conditions specified. To obtain a transform pair between  $x(\lambda)$  and  $C_n$ , the sum on the right-hand side must be reduced to one term or, at most, a finite number of terms independent of  $\lambda$ . This last is easily obtained since the integration depends only on the residues of  $Q_n(\lambda) \times Q_{-m-1}(\lambda)$  which are independent of  $\lambda$ . To obtain the former:

$$\frac{1}{2\pi j} \int_C Q_n(\lambda) Q_{-m-1}(\lambda) d\lambda = f(\delta_{n-m}) \quad (17)$$

must be satisfied where  $f(\delta_{n-m})$  is some realizable function or operator  $f$  of the Kronecker delta function  $\delta_{n-m}$ . In the

Taylor-Cauchy transform,  $f$  is the idem- function and it occurs because  $Q_n(\lambda)$  is a power function. Equation 17 is also satisfied when  $Q_n(\lambda)$  has orthogonal prop- erties. There are other properties of the moment function which will give rise to equation 17. Under these conditions equation 16 may be written as:

$$\sum_{n=0}^{\infty} C_n f(\delta_{n-m}) = \frac{1}{2\pi j} \int_C x(\lambda) Q_{-m-1}(\lambda) d\lambda \quad (18)$$

When  $f$  is a linear operation with respect to  $\delta_n$ :

$$\sum_{n=0}^{\infty} C_n f(\delta_{n-m}) = f \left[ \sum_{n=0}^{\infty} C_n \delta_{n-m} \right] = f(C_m) \quad (19)$$

whenever, in addition, the orders of sum- mation and  $f$  operation are interchange- able. The general transform is therefore given by:

$$x(\lambda) = \sum_{n=0}^{\infty} C_n Q_n(\lambda) \quad (20)$$

$$f(C_n) = \frac{1}{2\pi j} \int_C x(\lambda) Q_{-n-1}(\lambda) d\lambda \quad (21)$$

The linear operation  $f$  on  $C_n$ , namely  $f(C_n)$  and  $x(\lambda)$ , form a pair when properly treated and can be used in a manner similar to the Taylor-Cauchy transform for solving nonlinear differential equations describing physical nonlinear systems under the conditions specified in Appen- dix I. The functional form that the linear operator  $f$  takes is dependent upon the point of partition. The moment function  $Q_n(t)$  is also dependent upon the point of partition. The summary for  $k$ th-order real system is given in Table I.

It is to be noted that the phase parti- tioning at a given derivative is taken to mean that the given derivative and all higher-order derivatives remain in the left member while all other terms are trans- posed to the right.

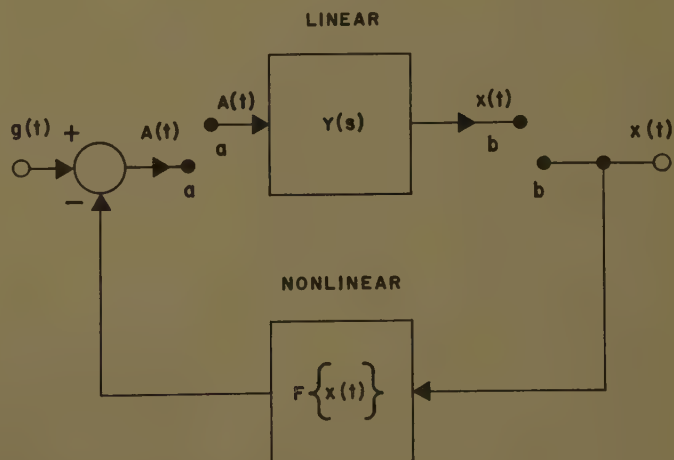


Fig. 3. The partitioned nonlinear system

EXAMPLE 1. PARTITIONING AN EQUATION

Given the nonlinear differential equation:

dx^2/dt^2 + dx/dt + bx + cx^2 = 1 (22)

1. Partitioning at the zero-th order derivative:

Qn(t) = 1/(gamma1 - gamma2) \* integral from 0 to t of tau^n [e^gamma1(t-tau) - e^gamma2(t-tau)] dtau (23)

The integral of equation 23 can be evaluated by making use of:

integral from 0 to t of tau^n e^-gamma tau dtau = (-1)^n d^n/dgamma^n [ (1 - e^-gamma t) / gamma ] (24)

where gamma1 and gamma2 are the roots of s^2 + as + b = 0.

2. Partitioning at the first derivative:

Qn(t) = t^(n+1)/(n+1) - integral from 0 to t of tau^n e^-a(t-tau) dtau (25)

3. Partitioning at the second derivative:

Qn(t) = t^(n+2)/(n+1)(n+2) (26)

EXAMPLE 2. APPLICATION OF TAYLOR-CAUCHY TRANSFORM

In this example apply the Taylor-Cauchy transform to case 3 of example 1. This gives:

wn = delta\_n - aw\_{n-1}/n - bw\_{n-2}/n(n-1) - C sum from k=2 to n-2 of wk wn-k-2 / ((k+1)(k+2)(n-k)(n-k-1)) (27)

and, by looking for the corresponding pairs in a table of Taylor-Cauchy transforms, one can obtain the time function x(t) corresponding to the response of equation 22.

Random Processes in Physical Systems

Consider a physical system such as shown in Fig. 2 and suppose the forcing function g(t) to be random with known statistical properties. The problem to be solved is to determine systematically the corresponding statistics of the response in terms of the statistics of the input. There are two theorems, proved elsewhere, for what follows which are useful in the general synthesis of random processes.

TRANSFORM-ENSEMBLE THEOREM

Given a random process g(t) with known statistical quantities and a deter-

ministic operator L such that if G(q) is the resulting transform-random process according to the relation

G(q) = L{g(t)} (28)

where q is the transform-random variable corresponding to the original random t, then

<G(q)>\_G = L[<g(t)>\_g] (29)

where <>\_G is the ensemble average with respect to G and <>\_g is the ensemble average with respect to g and commutative with respect to each other and the deterministic operator L, with the precaution that the subject of operation be the same for both operators.

The preceding theorem is useful for determining the first-order statistic of a system.

EXAMPLE 3. LINEAR BROWNIAN MOTION SYSTEM

Let v be the velocity of each of a large number of similar but independently acting iron particles subjected to an oscillating magnetic force of magnitude a and a random fluctuating force g(t) in a viscous medium giving rise to Doppler friction.

dv/dt + av = g(t) + a sin w\_0 t (30)

where a = f/m f = the coefficient of friction m = the mass of particle

Solution. In equation 29 let the operator L denote the Laplace operation. In a Brownian motion phenomenon the particles are all assumed to have started with the same initial velocity, v\_0. Thus:

G(s) + a w\_0 / (s^2 + w\_0^2) + V\_0 / (s + a) (31)

Utilizing the inverse of equation 29 for this case

<v(t)>\_v = L^-1[<V(s)>\_v] (32)

The result

<v(t)>\_v = L^-1 [ <G(s)>\_G / (s + a) + a w\_0 / ((s + a)(s^2 + w\_0^2)) + v\_0 / (s + a) ]

If g(t) is stationary:

<G(s)>\_G = <g(t)>\_g / s

<v(t)>\_v = <g(t)>\_g [ 1 - e^-at ] + e^-at [ sin w\_0 t - w\_0 cos w\_0 t ] / (a^2 + w\_0^2) + v\_0 e^-at

Higher-order statistics can easily be obtained by appealing to the theorems given below. In order to discuss the theorem, first denote by c^k{ } the convolution transform of the kth order in the complex s plane this is:

c{F(s)} = F(s) \* F(s) c^2{F(s)} = F(s) \* F(s) \* F(s) ... etc.

where \* denotes the convolution in the complex plane.

Higher-order statistic transform theorem. Given a random process g(t) and its transform G(q), the mean nth ensemble average of the given process is:

<g^n(t)>\_g = L^-1[<c^{n-1}{G(q)}>\_G]

The proof of this theorem is obtained by induction from the previous theorem with the aid of the convolution transform.

EXAMPLE 4. SECOND STATISTIC SIMPLE STOCHASTIC PROCESS

Consider a simple stochastic process defined by the differential equation initially at rest:

dx/dt = g(t) e^-at; a > 0

where g(t) is an ergodic Gaussian random process with zero mean and unit variance. The problem is to obtain the variance of the response when the power spectrum is uniform over the frequency spectrum.

Solution. The Laplace transform of equation 39 for x(0) = 0 is:

X(s) = G(s + a) / s

Taking the convolution transform

Table I

Partition Point	Moment Function
At kth derivative	power function
At kth + (k-1) st derivatives	nonperiodic function of exponential type
At kth + (k-2) nd derivatives	periodic function of exponential type
At all linear derivatives	combinations and linear combinations of power functions, and nonperiodic functions of exponential type



both sides of 40 and applying the foregoing theorem gives

$$\begin{aligned} \langle x^2(t) \rangle_x &= L^{-1}[\langle \mathcal{C}X(s) \rangle_x] \\ &= L^{-1}\left[\frac{\langle G(a)G(s+a) \rangle_g}{s} \right] \quad (41) \end{aligned}$$

To evaluate equation 41, use is made of the formula:<sup>10</sup>

$$\begin{aligned} \langle G(z+a)G(s-z+a) \rangle_g &= \\ L_z L_s \{ \langle g(t)g(t-\tau) \rangle_g e^{-a\tau} e^{-2a\tau} \} \quad (42) \end{aligned}$$

where the complex variable  $s$  corresponds to  $t$ , a real variable, and  $z$  corresponds to  $\tau$ , another real variable. Noting that the power-density spectrum is flat, i.e., white noise (equation 42) reduces to the simple result:

$$\langle G(z+a)G(s-z+a) \rangle_g = \frac{1}{s+2a} \quad (43)$$

Putting this into 41 gives:

$$\langle x^2(t) \rangle_x = L^{-1}\left[\frac{1}{s(s+2a)}\right] \quad (44)$$

or

$$\langle x^2(t) \rangle_x = \frac{1}{2a} [1 - e^{-2at}] \quad (45)$$

Since the mean value of the process is zero, the variance  $\sigma^2$  is:

$$\sigma^2 = \left[ \frac{1}{2a} (1 - e^{-2at}) \right] \quad (46)$$

## Relations between Wiener-Lee-Bose Theory and Wolf-Ku Theory

The Wiener theory of characterizing a nonlinear system consists, briefly, of obtaining a set of coefficients which are capable of physical measurement. These coefficients are obtained by subjecting the system under test to a white-noise probe. Because of the properties of white noise these coefficients determine uniquely the transmission characteristics of the given system. Since many systems with different physical configurations can exhibit the same transmission characteristics, the coefficients therefore do not uniquely determine the physical configuration of the system. The given system can, however, be found uniquely by a variational procedure which will be described.

Consider Fig. 4. This scheme was developed by Wiener, Lee, and Bose to test a nonlinear system so that information obtained from the response could be used to synthesize a new system having the same response to white noise as the original system. Since the power-density spectrum of white noise is flat over the infinite frequency interval, the system is effectively being exposed to every physical

signal. It is well known that, in order to know the response of a nonlinear system to every input, it must be tested with every input.<sup>2</sup> This is unlike the behavior of a linear system that requires only an impulse as an input to determine its response to every input. This last occurs because of superposition. Hence, if the response to an arbitrary input is required, one would use the well-known convolution integral with the impulse response as a weighting function.

The Wiener theory fails to give back the original system configuration since the white-noise probe does not consider the effect of initial conditions. In linear systems the response is essentially independent of the initial boundary conditions. In nonlinear systems a change in initial conditions produces a change in the response. The most general test of nonlinear systems would be, then, to subject the input to every conceivable input and every conceivable initial value.

Fortunately, however, the use of white noise and parameter variation produces the same effect by generating a set of coefficients which not only uniquely defines the transmission, but also the con-

figuration. Now consider the linear case.

## LINEAR CASE

Let  $Y(s)$  be the transfer function of the physical system. Thus:

$$Z(D)x(t) = g(t) \quad (47)$$

describes the dynamic behavior of the linear system. Hence, by inversion:

$$x(t) = \int_{-\infty}^{\infty} g(t-\tau)y(\tau)d\tau \quad (48)$$

Similarly, from Fig. 4:

$$v_n(t) = \int_{-\infty}^{\infty} g(t-\gamma)h_n(\gamma)d\gamma \quad (49)$$

Multiplying equations 48 and 49 and averaging:

$$\begin{aligned} \overline{x(t)v_n(t)} &= \\ \int_0^{\infty} \int_0^{\infty} \overline{g(t-\tau)g(t-\gamma)y(\tau)h_n(\gamma)} d\gamma d\tau \quad (50) \end{aligned}$$

Noting that white noise is a Gaussian process in this case, the outputs of the orthogonal filter are also Gaussian because

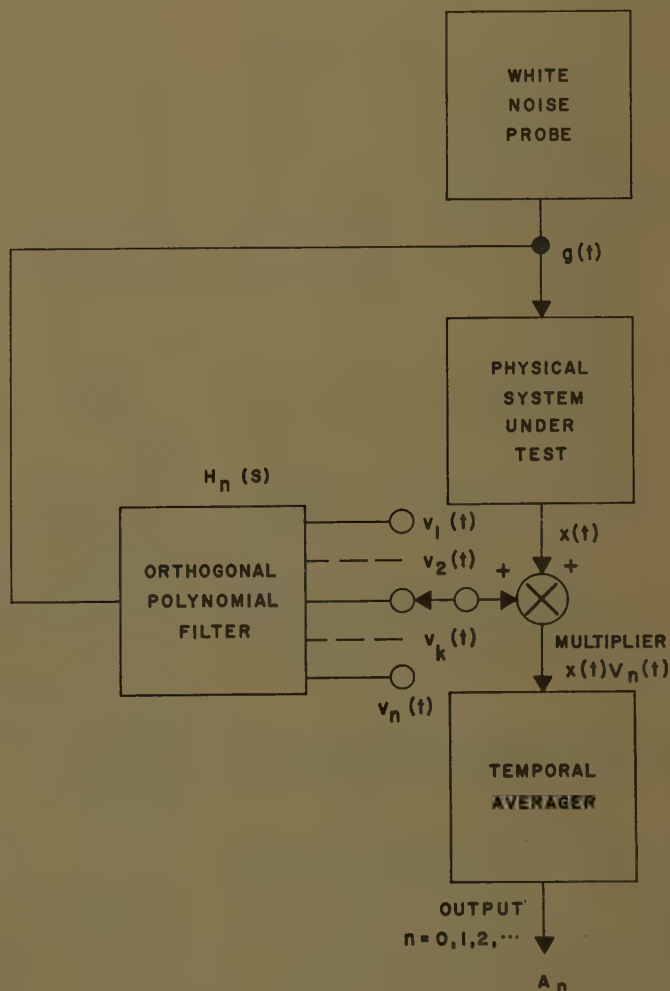


Fig. 4. The Wiener-Lee-Bose system

of linearity. Furthermore the outputs are statistically independent and all have the same variance. This obtains from the orthogonality and the nature of the network's nondissipation. Because of ergodicity, replace time averages with ensemble averages. Therefore equation 50 becomes

$$\overline{x(t)v_n(t)} = \int_0^\infty \int_0^\infty \overline{g(t-\tau) \times g(t-\gamma)} y(\tau) h_n(\gamma) d\tau d\gamma \quad (51)$$

The average inside the integral is recognized as the autocorrelation function. For white noise this is:

$$\overline{g(-\tau)g(t-\gamma)} = \delta(\tau-\gamma) \quad (52)$$

where  $\delta(\tau-\gamma)$  denotes the Dirac delta function with power density.

Substituting equation 52 into equation 50 and simplifying gives:

$$\overline{x(t)v_n(t)} = \int_0^\infty y(\tau) h_n(\tau) d\tau \quad (53)$$

The left-hand side is obtained by the Wiener theory:

$$\overline{x(t)v_n(t)} = A_n \quad (54)$$

To display the dependence of  $y$  on its parameters:

$$y(\tau) = y(\tau; \alpha_1, \dots, \alpha_k; \beta_1, \dots, \beta_m) \quad (55)$$

where  $\alpha_j$  and  $\beta_i$  are the time constants of the linear system. It is then evident that if  $n = k + m$  measurements were made, the  $\alpha_j$  and  $\beta_i$  could be uniquely determined from equation 56 by simultaneous solution of the set of algebraic equations which result.

$$A_n = \int_0^\infty y(\tau; \alpha_1, \dots, \alpha_k; \beta_1, \dots, \beta_m) \times h_n(\gamma) d\tau \quad (56)$$

To obtain the values of  $\alpha_j$  and  $\beta_i$  in terms of the actual given components, use equation 57:

$$\delta A_n = \int_0^\infty \delta y(\tau; \alpha_1, \dots, \alpha_k; \beta_1, \dots, \beta_m) \times h_n(\tau) d\tau \quad (57)$$

where the variation  $\delta$  is made with respect to the proper number of time constants according to the total number of components in the system.

## Physical Interpretation of Integral Equations of Identification

Equations 56 and 64 are called the integral equations of identification since, as pointed out, the system parameters for linear and nonlinear systems can be determined from them by measurement of only the  $A_n$  statistics.

Examination of these equations shows that they are amenable to the solution of the system function  $y(\tau)$ , in the linear case, in terms of the orthogonal functions of  $h_n(\tau)$  as linear combinations of the latter. That is, solving for  $y(\tau)$  in equation 53 gives:

$$y(\tau) = \sum_{n=0}^\infty A_n h_n(\tau) \quad (53A)$$

subject only to the constraint given by

$$\int_0^\infty h_k(\tau) h_n(\tau) d\tau = \delta_{k-n} \quad (53B)$$

where, as before,  $\delta_{k-n}$  is a Kronecker delta function.

It is therefore clear that equation 53(A) denotes a method of measuring the impulse response in terms of the impulse responses of the orthogonal filters by measuring only the statistics associated with the system under test, as depicted in Fig. 4, when the random source is stationary and white. Suppose now that  $g(t)$  is an arbitrary random source so that its correlation function is

$$\phi(\tau-\gamma) = \overline{g(t-\tau)g(t-\gamma)} \quad (53C)$$

In terms of equation 53(C), equation 50 becomes

$$A_n = \int_0^\infty \int_0^\infty \phi(\tau-\gamma) y(\tau) h_n(\gamma) d\tau d\gamma \quad (53D)$$

If the system under test is fixed to, say, a straight-through connection, then:

$$y(\tau) = \delta(\tau) \quad (53E)$$

Thus:

$$A_n = \int_0^\infty \phi(\gamma) h_n(\gamma) d\gamma \quad (53F)$$

from which

$$\phi(\gamma) = \sum_{n=0}^\infty A_n h_n(\gamma) \quad (53G)$$

This gives a convenient method of measuring the autocorrelation function in a manner similar to system testing given in equation 53(A). Such a correlator has been tested experimentally and its results will be reported elsewhere.

Therefore, the integral equations of identification give rise to convenient methods of obtaining system attenuation and phase characteristics and correlation functions depending on whether the noise source is fixed or the system under test is fixed. Similar interpretation exists for the nonlinear case.

## THE NONLINEAR CASE

Consider equation 4 written as

$$x(t) = \int_0^\infty a(t-\tau) y(\tau) d\tau \quad (58)$$

where  $a(t)$  is the modified forcing function. Let  $a(t)$  be expanded into a series of complete orthogonal functions

$$a(t) = \sum_{k=0}^\infty a_k P_k(t) \quad (59)$$

where

$$\int_0^\infty P_k(t) P_m(t) dt = \delta_{k-m} \quad (60)$$

and

$$a_k = \int_0^\infty a(t) P_k(t) dt \quad (61)$$

Expansion 59 is justified since  $P_k(t)$  is taken analytic for all  $k$ . Using equation 49 and multiplying and averaging:

$$\overline{x(t)v_n(t)} = \int_0^\infty \int_0^\infty \overline{a(t-\tau)g(t-\gamma)} \times y(\tau) h_n(\gamma) d\tau d\gamma \quad (62)$$

Now let

$$\overline{a(t-\tau)g(t-\gamma)} = \phi(\tau-\gamma) \quad (63)$$

so that

$$\overline{x(t)v_n(t)} = \int_0^\infty \int_0^\infty \phi(\tau-\gamma) y(\tau) h_n(\gamma) d\tau d\gamma \quad (64)$$

Solution of equation 64 requires a knowledge of the expansion of  $\phi(\tau-\gamma)$ . Consider one case patterned after the theory given in references 2 and 3. Let the right side of equation 63,  $\phi$ , occupy the same place, with respect to a random process, as  $a(t)$  occupies with respect to a deterministic process. Then, under certain conditions,<sup>2-3</sup>  $\phi(\tau)$  can be expanded as follows:

$$\phi(\tau) = \sum_{n=0}^\infty C_n \tau^n \quad (65)$$

Noting that<sup>21</sup>

$$\phi(\tau-\gamma) = e^{-\gamma D} \phi(\tau) \quad (66)$$

where  $D$  is differential time operator  $d/d\tau$  and substituting the right side of equation 65 into 66 for  $\phi(\tau)$ , results in:

$$\phi(\tau-\gamma) = e^{-\gamma D} \sum_{n=0}^\infty C_n \tau^n \quad (67)$$

Putting equation 67 into 64 gives:

$$\overline{x(t)v_n(t)} = \int_0^\infty \int_0^\infty \left( e^{-\gamma D} \sum_{n=0}^\infty C_n \tau^n \right) \times y(\tau) h_n(\gamma) d\tau d\gamma \quad (68)$$

Interchanging the order of summation and integration yields

$$\overline{x(t)v_n(t)} = \sum_{n=0}^\infty C_n \int_0^\infty \int_0^\infty (e^{-\gamma D} \tau^n) y(\tau) \times h_n(\gamma) d\tau d\gamma \quad (69)$$



is now convenient to define a new parameter,  $M_n(\gamma)$ , as:

$$M_n(\gamma) = \int_0^\infty (e^{-\gamma \tau} \tau^n) y(\tau) d\tau \quad (70)$$

$$M_n(\gamma) = \int_0^\infty \tau^n y(\gamma + \tau) d\tau \quad (70A)$$

This parameter represents the modified moments of  $y(\tau)$ , the impulse response of the linear part of the system for a random process, in the same manner that  $Q_n(t)$  represents the moments of equation 9 for a deterministic process. Using equation 69 with 69 yields

$$M_n(t) = \sum_{n=0}^\infty C_n K_n \quad (71)$$

where  $K_n$  is defined by equation 72:

$$K_n = \int_0^\infty M(\gamma) h_n(\gamma) d\gamma \quad (72)$$

The  $C_n$ 's are obtained recursively following a similar procedure to one given previously.<sup>2-4</sup> Note that  $x(t)$  can be expanded into a power series<sup>9</sup> such as:

$$x(t) = \sum_{n=0}^\infty f_n t^n \quad (73)$$

When the random process is a Brownian motion type. Then, noting that equation 59 can also be expanded into a power series:

$$y(t) = \sum_{n=0}^\infty b_n t^n \quad (74)$$

where  $b_n$  are the Taylor coefficients of the right side of equation 59 namely:

$$b_n = \frac{1}{n!} \frac{d^n}{dt^n} \left[ \sum_{k=0}^\infty a_k P_k(t) \right]_{t=0} \quad (75)$$

The last coefficients may also be obtained recursively. Rewriting equation 58 as

$$y(t) = \int_0^\infty a(\tau) y(t-\tau) d\tau \quad (76)$$

and using equation 74 yields:

$$y(t) = \sum_{n=0}^\infty b_n Q_n(t) \quad (77)$$

When the partition is made at the highest-order derivative,  $Q_n(t)$  is a power function.<sup>2-3</sup> In this case, recurrence relations are obtained relating the  $b_n$  to the  $a_n$  given in equation 59. The  $f_n$  are automatically obtained from equation 77. To relate the  $C_n$  to the  $b_n$ , appeal to equation 63. The latter yields the appropriate recurrence relation to complete the solution. The recurrence relations will not be dealt with in detail, however, a great deal of simplification can be effected in the form

of the recurrence relations if  $a(t)$  is assumed to be a power series at the start. In this case the Taylor-Cauchy transform can be applied and will now be considered with the detection of faults in complex physical systems.

## General Synthesis Problem

In this section, a theory leading to a general synthesis procedure of both linear and nonlinear systems will be considered, the problem being that of uniqueness.

In the problem of analysis the solution is always unique for linear systems and for a certain class of nonlinear systems which are physically realizable.<sup>2-3</sup> In the problem of synthesis the solution is not generally unique. The reason for this lack of uniqueness seems to be related to a corresponding lack of information concerning such things as initial conditions, forcing functions, distribution of currents flowing in the interior of the structure and the number of components comprising the configuration. Generally, in linear synthesis procedures, the forcing function is tacitly given when the procedure involves poles and zeros in the complex frequency plane, otherwise none of the other information is usually given. This allows for an ambiguity in construction of the configuration giving rise to as many as an infinite variety of solutions.

In the previous section it was pointed out that the Wiener theory by itself led to a synthesis procedure which was not unique. The Wiener structure combined Laguerre and Hermite networks with multipliers and adders to produce a response to white noise which was identical to the system under test. By combining the Wiener-Lee-Bose theory with the Wolf-Ku theory a set of measured values are related to the transmission coefficients of the structure which are in turn related to the time constants of the structure. However, this knowledge does not give a unique structure of the system under test since once again many systems can have equivalent transmission with such time constants. To define the nature of these

time constants in order to define uniquely the structure, the system should be tested in the following way. Expose the system to a white-noise probe as a forcing function and, simultaneously, vary the boundary conditions according to another white-noise probe generated independently; see Fig. 5. The unique synthesis is obtained by a statistical analysis of the output. This might be achieved by subjecting the output to an analysis similar to the Wiener scheme with the analysis repeated with orthogonal components obtained from both white-noise probes. This problem has not yet been satisfactorily solved. The main difficulty stems from the introduction of the second noise probe.

As remarked earlier, if the structure configuration is given a priori, the unique synthesis problem is reduced to a procedure involving two sets of measurements of the Wiener type; the second is obtained by repeating the first after a variation has been performed on the system in a specified manner.

## Synthesis to Special Waveforms-Pole-Zero Method

In the previous section the synthesis of networks and systems under the influence of a white-noise probe was discussed. Now, a synthesis procedure is indicated for nonlinear systems given a specified waveform as an input, boundary values, a distribution of singularities in the complex plane, and the transmission characteristics of the linear part. The procedure is based on a theory outlined for stable systems in reference 2 and continued in reference 7, however, this procedure has a much simpler concept.

Let  $Y(s)$  be the transfer function of the linear part and let  $\{s_n\}$  be a set of singularities which are generally supposed to be finite in number; then, given a specified forcing function  $g(t)$ , what is the nature of the nonlinear part  $F$  to satisfy these conditions?

The synthesis procedure under these conditions is:

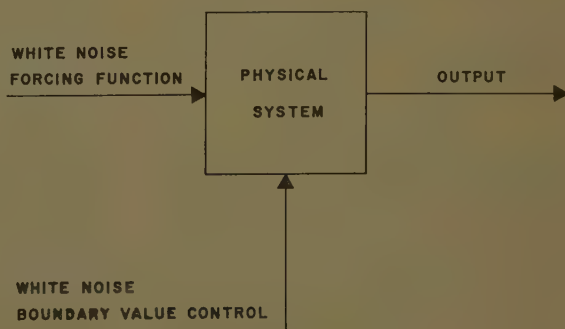


Fig. 5. A procedure for uniquely determining the synthesis of a given structure

1. From the singularities and initial conditions determine the coefficients  $C_n$  in closed form using reference 2.
2. Utilizing the  $C_n$  coefficients, determine the response of the system using the moment theorem.
3. The nonlinear part of the system is then obtained using the canonical equation 1.
4. The canonical system is then obtained easily by reference to Fig. 2.

#### EXAMPLE

Let the linear plant of a system be given as

$$Y(s) = \frac{1}{s+2} \quad (78)$$

If the system is forced by

$$g(t) = e^{-2t} \quad (79)$$

find the coefficients of the nonlinear terms so that the response has a pole at  $s = -1$  and the nonlinear part has no derivatives and does not exceed degree 2.

The solution can be effected by using the recurrence equations of references 2-4, or it may be solved directly, in which case the solution is almost trivial. For instance, from

$$Z(S) = \frac{1}{Y(S)} = S+2 \quad (80)$$

the canonical form is obtained:

$$\frac{dx}{dt} = 2x + F(x) = g(t) \quad (81)$$

or

$$F(x) = g(t) - \dot{x} - 2x \quad (82)$$

Noting

$$F(x) = \sum_{k=0}^n a_k x^k$$

then

$$\sum_{k=0}^2 a_k x^k = g(t) - \dot{x} - 2x \quad (83)$$

From the problem:

$$X(s) = \frac{1}{s+1} \quad (84)$$

from which

$$x(t) = e^{-t} \quad (85)$$

so that

$$a_0 + a_1 e^{-t} + a_2 e^{-2t} = e^{-2t} - e^{-t} \quad (86)$$

Since equation 86 is an identity, it is clear that:

$$a_0 = 0 \quad (87)$$

$$a_1 = -1 \quad (88)$$

$$a_2 = 1 \quad (89)$$

Thus the nonlinear component satisfying the problem requirements is:

$$F(x) = -x + x^2 \quad (90)$$

#### Conclusions

This paper has presented some aspects of some of the recent advances in the analysis and synthesis of nonlinear systems. The concepts involved in the Wiener theory and those developed by the author have been reviewed. It is evident that the two theories have different approaches looking toward common goals. On one hand, the Wiener theory considers the synthesis problem in a statistical form, while on the other hand the author's theory considers the concept of partition which is a deterministic process. However, the two theories may be combined to produce a new theory for determining the structural configuration internal to a given system. This theory therefore finds application in the fault detection and diagnosis problem in complex systems. In addition to its practical value the new theory relates coefficients capable of being measured to coefficients capable of being calculated, which gives a means of mechanizing new measuring systems with high accuracy.

By interpreting the solutions of the integral equations of identification, one can devise highly accurate measuring systems to determine attenuation and phase characteristics for linear systems, and correlation functions of stationary random functions. This leads to the concept of "almost instantaneous" and "instantaneous" correlators.

A complete synthesis theory can be developed using the canonical system as a basis. From here it is easy to show that the response of this system is unaltered if the linear and nonlinear boxes are interchanged providing each operator is also inverted. It is possible to show that by making the appropriate partitions almost any desired configuration can be achieved in both single- and multiple-loop feedback systems.

It is interesting to note that the absolute test for stability of a nonlinear system consists of utilizing a double white-noise probe (Fig. 5) and examining the position of the resulting singularities in the complex frequency plane as described in reference 7.

Some very important problems are still left unsolved; for instance, the solution of the recurrence equations in closed form, development of other transforms like the Taylor-Cauchy to give elementary and higher transcendental types of solutions for nonlinear systems, and using

more efficiently Kronecker's theorem (developed in reference 2) as a tool in the synthesis of nonlinear systems.

It is possible to extend the synthesis problem, under suitable restrictions, to the design of adaptive control systems.

#### Appendix I. Restrictions on System

1. The linear operator  $Z(D)$  has an impulsive response  $y(t)$  of exponential type and order 1.
2. The forcing function  $g(t)$  is a function of exponential type and order 1.
3. The nonlinear function  $F(x, \dot{x}, \dots)$  is single-valued, continuous, and satisfies the following conditions.

Given a pair of positive constants,  $M$  and  $a$ , and two continuous functions,  $v(t)$  and  $u(t)$ , which are asymptotically like  $e^{-at}$  then  $F$  must satisfy the condition that

$$|F\{v(t)\} - F\{u(t)\}| < Me^{-at}|v(t) - u(t)|$$

for all  $t > 0$ .

#### Appendix II. Partitioning to Obtain Power-Series Solutions

In the partition theory of solving ordinary nonlinear differential equations of a certain class, all the nonlinear terms are transposed to the right side. The solution is a linear combination of the moments of the folded impulse response of the linear terms in the open interval  $(0, t)$ . It is noted that if all but the highest-order derivative is transposed to the right member, power-series solution results because the resulting moment function is a power function. It is possible to obtain power-series solutions by transposing all terms except any one linear term to the right side. However, in the case of linear terms of order less than the highest-order derivative, it is necessary to reorder the indexes prior to obtaining the recurrence equation. If prior reordering of indexes is not according to the order of the term remaining in the partition, the resulting recurrence relation will be backward.

The partitioning of an equation at the highest derivative has the benefit of giving the desired coefficients automatically as a function of the lower-order ones without prior reordering of indexes.

#### References

1. THEORY OF NONLINEAR CONTROL, Y. H. KU, "Proceedings, First International Symposium on Automatic Control," Moscow, Aug. 1960, Butterworth's Scientific Publications, London, England, 1961; also *Journal, Franklin Institute, Philadelphia, Pa.*, vol. 271, Feb. 1961, pp. 108-44.
2. A MATHEMATICAL THEORY FOR THE ANALYSIS OF A CLASS OF NONLINEAR SYSTEMS, Alfred A. Wolf, *Doctoral Dissertation*, The University of Pennsylvania, Philadelphia, Pa., June 1958.
3. RECURRENCE RELATIONS IN THE SOLUTION OF A CERTAIN CLASS OF NONLINEAR SYSTEMS, Alfred A. Wolf, *AIEE Transactions*, pt I (Communication and Electronics), vol. 78, 1959 (Jan. 1960 section), pp. 830-34.



GENERALIZED RECURRENCE RELATIONS IN THE ANALYSIS OF NONLINEAR SYSTEMS, Alfred A. Wolf. *Ibid.*, vol. 80, Sept. 1961, pp. 383-87.

ANALYSIS OF TRANSCENDENTAL NONLINEAR SYSTEMS, Alfred A. Wolf. *Ibid.*, vol. 79, Nov. 1960, pp. 449-51.

ON THE SIGNIFICANCE OF TRANSIENTS AND STEADY STATE BEHAVIOR IN NONLINEAR SYSTEMS, A. Wolf. *Proceedings*, Institute of Radio Engineers, New York, N. Y., vol. 47, no. 10, Oct. 1959, pp. 1785-86.

A STABILITY CRITERION FOR NONLINEAR SYSTEMS, Y. H. Ku, A. A. Wolf. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 78, July 1960, pp. 144-48.

TAYLOR-CAUCHY TRANSFORMS FOR ANALYSIS OF A CLASS OF NONLINEAR SYSTEMS, Y. H. Ku, A. Wolf, J. H. Dietz. *National Convention Record*, Institute of Radio Engineers, New York, N. Y., pt. 2 (*Circuit Theory*), 1959, pp. 49-61; also *Proceedings*, Institute of Radio Engineers, vol. 48, no. 5, May 1960, pp. 912-22.

TRANSFORM-ENSEMBLE METHOD FOR THE

ANALYSIS OF LINEAR AND NONLINEAR SYSTEMS WITH RANDOM INPUTS, Y. H. Ku, A. A. Wolf. *Proceedings*, National Electronics Conference, Chicago, Ill., vol. 15, 1959, pp. 441-45.

10. ON POLES AND ZEROS OF A RANDOM PROCESS IN LINEAR AND NONLINEAR SYSTEMS, A. A. Wolf. *Ibid.*, vol. 16, 1960, pp. 268-78.

11. A THEORY OF NONLINEAR SYSTEMS, A. G. Bose. *Technical Report no. 309*, Research Laboratory of Electronics, Massachusetts Institute of Technology, Cambridge, Mass., May 15, 1960.

12. NONLINEAR PROBLEMS IN RANDOM THEORY (book), N. Wiener. John Wiley & Sons, Inc., New York, N. Y., 1958.

13. CYBERNETICS (book), N. Wiener. John Wiley & Sons, Inc., 1948.

14. APPLICATION OF STATISTICAL METHODS TO COMMUNICATION PROBLEMS, Y. W. Lee. *Technical Report no. 181*, Research Laboratory of Electronics, Massachusetts Institute of Technology, Sept. 1, 1950.

15. THEORY OF NONLINEAR TRANSDUCERS, H. E.

Singleton. *Technical Report no. 160*, Research Laboratory of Electronics, Massachusetts Institute of Technology, Aug. 12, 1950.

16. THEORIE UND ANWENDUNG DER LAPLACE TRANSFORMATION (book), G. Doetsche. Dover Publications, Inc., New York, N. Y., 1943.

17. ENTIRE FUNCTIONS (book), R. P. Boas. Academic Press, New York, N. Y., 1954.

18. ON A SYSTEMATIC APPROXIMATION TO THE PARTITION METHOD FOR ANALYSIS OF A CLASS OF NONLINEAR SYSTEMS, Y. H. Ku, A. A. Wolf, J. H. Dietz. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 79, July 1960, pp. 183-91.

19. ON THE THEORY OF THE BROWNIAN MOTION, G. E. Uhlenbeck, L. S. Ornstein. *Physical Review*, New York, N. Y., Sept. 1, 1930, pp. 823-41.

20. THE FOURIER INTEGRAL AND CERTAIN OF ITS APPLICATIONS (book), N. Wiener. Dover Publications, Inc., 1933.

21. THE CONVOLUTION TRANSFORM (book), D. V. Widder, I. I. Hirschman. Princeton University Press, Princeton, N. J., 1955.

## Discussion

Y. H. Ku (Moore School of Electrical Engineering, University of Pennsylvania, Philadelphia, Pa.): In connection with the partition theory discussed in the beginning of Dr. Wolf's paper, I did suggest the name "auxiliary forcing function" for  $A(t)$  used in equations 4 and 5 and shown in the block diagram, Fig. 2, of the paper. I wish to point out that in the nonlinear differential equation 1, transposition is allowable for a class of problems where the solution is unique and analytic, although superposition would certainly fail to apply. In a recent paper<sup>1</sup> I have discussed this point in connection with the feedback system block diagram. Let the nonlinear differential equation be

$$+k_1c' + k_2c = r - N(c', c) = r - b = e \quad (91)$$

where  $c$  denotes the output signal,  $c'$  and  $c''$  denote the first and second derivatives of  $c$  with respect to time, and  $N(c', c)$  denotes a nonlinear function of  $c'$  and  $c$ . The symbols  $r$ ,  $b$ , and  $e$  correspond to the input, feedback signal, and error or actuating signal respectively. Note that the actuating signal  $e(t)$  corresponds to a particular choice of the auxiliary forcing function  $A(t)$  mentioned by Dr. Wolf's paper.

We can rewrite equation 91 in two different ways:

$$+k_1c' = r - k_2c - N(c', c) = r - b_1 = e_1 \quad (92)$$

$$= r - k_1c' - k_2c - N(c', c) = r - b_2 = e_2 \quad (93)$$

Thus,  $e_1$  and  $e_2$  correspond to  $A_1(t)$  and  $A_2(t)$ , the two other modified auxiliary functions, respectively. Suitable block diagrams can be drawn corresponding to equations 2 and 3. So the method is not limited to partitioning at the highest derivative alone, though the highest derivative should be included on the left-hand side of equations 91-93. Note that the original nonlinear differential equation would be

$$+k_1c' + k_2c + N(c', c) = r \quad (94)$$

which may represent an open-loop system. The advantage of the Taylor-Cauchy transform method is in its possibility of many applications. For instance, it can be used with the Transform-Ensemble method

mentioned in reference 9 of the paper. It can be readily used for the solution of varying-parameter systems.<sup>2</sup> The writer has great hope in Dr. Wolf's extension of the above-mentioned approach to the synthesis of nonlinear systems. Maybe it is in the field of synthesis rather than the field of analysis that the new method will be most useful.

## REFERENCES

1. ON NONLINEAR NETWORKS WITH RANDOM INPUTS, Y. H. Ku. *Transactions*, Institute of Radio Engineers, New York, N. Y., vol. CT-7, no. 4, Dec. 1960, pp. 479-90.

2. TAYLOR-CAUCHY TRANSFORMS FOR ANALYSIS OF VARYING-PARAMETER SYSTEMS, Y. H. Ku. *Proceedings*, Institute of Radio Engineers, vol. 49, no. 6, June 1961, pp. 1096-97.

R. L. Cosgriff (Ohio State University, Columbus, Ohio): The author is to be congratulated for focusing attention upon series solutions of nonlinear differential equations. This discussion is concerned with the real time series solutions and the comparison of this technique with the Taylor-Cauchy transform method.

Nonlinear differential equations with deterministic driving functions can be solved as accurately as desired providing one is willing to expend the energy and a solution exists. Many methods can be employed, and a judicious choice of method for a given problem can greatly reduce the labor involved. One of the classical methods for the solution of linear differential equations is the method of Frobenius, that gives the solution in terms of a Taylor series. For a complete discussion see reference 1. This method when extended to the nonlinear<sup>2</sup> case yields results identical to those obtained by the Taylor-Cauchy transforms. Disregarding terminology, the basic steps of this transform technique and those involved in the method of Frobenius are identical. In both cases a Taylor expansion of all variables involved in the differential equation must exist. (The collective coefficients of the series for each of the various variables is the transform of that variable.) Equating all coefficients, known and unknown, of like powers of  $t$  arising after the substitution of time series into the differential equation corresponds to equating the transforms of the variables involved. Determining the transform of the unknown variable corre-

sponds to determining the coefficients of the time series of this unknown variable.

Consider the author's example

$$\frac{d^2x}{dt^2} + a \frac{dx}{dt} + bx + cx^3 - 1 = 0$$

$$x = \dot{x} = 0 \text{ at } t = 0$$

Let

$$x = \sum_{n=2}^{\infty} w_n t^n$$

$$w_0 = w_1 = 0$$

the series for  $x^2$  becomes

$$x^2 = \sum_{n=2}^{\infty} \sum_{m=2}^{n-2} w_m w_{n-m} t^n$$

Upon substitution of the assumed series one has

$$\sum \left[ n(n-1)w_n + a(n-1)w_{n-1} + bw_{n-2} + c \sum_{m=2}^{n-4} w_m w_{n-m-2} \right] t^{n-2} - 1 = 0$$

Thus, the coefficients for the unknown  $x$  (the transform of  $x$ ) are given by

$$w_n = -\frac{aw_{n-1}}{n} - b \frac{w_{n-2}}{n(n-1)} - c \sum_{m=2}^{n-4} \frac{w_m w_{n-m-2}}{n(n-1)}$$

$$\text{for } n > 2$$

The above recursive relationship gives the unknown coefficients (the transform) of  $x$  and is identical to equation 27 once the  $w$ 's of equation 27 are properly defined or the expression is corrected. For  $a=b=c=1$ ,  $w_2=1/2$ ,  $w_3=-1/16$ ,  $w_4=0$ ,  $w_5=1/120$ ,  $w_6=-1/720-1/120=-7/720$ , therefore

$$x = t/2^2 - 1/16t^3 + 1/120t^5 - 7/720t^6 + \dots$$

Notice the two terms of  $w_6$ . The first would occur if  $c=0$  (the linear differential equation) and the second is due to the nonlinear term. Observe the large magnitude due to the nonlinear term. For  $n>6$  the nonlinear terms predominately determined the values of the coefficients of the differential equation. Thus, it can be concluded that the rate of convergence must be considered,

and one can easily demonstrate that this is far more of a limitation in the solution of nonlinear problems than it is in the case of linear problems.

Returning to the author's development it is apparent that the special techniques such as partition theory, moment theory, and the Taylor-Cauchy transform can be bypassed by using the method of Frobenius.

The author considers desirable requirements for other transforms equation 16 through 21. These conditions describe functional groups convenient for the solution of differential equations. In the real domain desirable requirements for functional groups can be described; for example,

$$\frac{df_n}{dt} = \sum_{M_1}^{M_2} k_m f_m \quad M_1 \leq M_2 \leq n$$

$$f_m f_n = \sum_{\alpha}^{\beta} K_{\alpha} f_{\alpha}$$

with the requirement  $m \leq n \leq \alpha \leq \beta$  where  $k$ 's and  $K$ 's are constants. If these conditions are met and the driving function,  $y$ , can be expanded in terms of the summation

$$y = \sum c_{\alpha n} f_n(t)$$

the solution for  $x$  will take the form

$$x = \sum c_{\beta n} f_n(t)$$

and the recursive relationships for the  $c_{\beta n}$  will exist. Simple functions are

$$f_n = t^n$$

$$f_n = e^{-\alpha n t} \quad (\alpha \text{ is a constant})$$

$$f_n = t^n e^{-\alpha t}$$

I wonder whether these real time requirements are the same as the requirements based upon transform theory.

The importance of Wolf's work with transform methods is that it introduces a new facet concerning the solution of nonlinear differential equations. As each facet is exploited new insights and techniques develop which expand our knowledge of nonlinear systems.

# REFERENCES

1. ORDINARY DIFFERENTIAL EQUATIONS, (book), E. L. Ince. Dover Publications, Inc., New York, N. Y., pp. 396-403.
2. NONLINEAR CONTROL SYSTEMS (book), R. L. Cosgriff. McGraw-Hill Book Company, Inc., New York, N. Y., 1958, pp. 15-16.

Louis F. Kazda and A. Y. Bilal (University of Michigan, Ann Arbor, Mich.): Dr. Wolf has called the attention of control system engineers to the partitioning method and its application to nonlinear differential equations. It is a systematic way of solving nonlinear differential equations, and lends itself to digital computation. Probably the greatest contribution of the method is that it permits an insight into the synthesis of nonlinear systems.

A review of the author's published work, however, would reveal that certain questions have been left unanswered, and which, the novice, trying to utilize the presented material, would naturally like to have summarized. These questions are:

1. What insight can be gained about the stability of a system for a general class of forcing functions utilizing the partitioning method of analysis?
2. Under what conditions would partitioning at points other than the highest derivative be desirable?
3. Consider for example the nonlinear differential equation

$$\ddot{x} + \omega_0 x + x^3 = 0$$

which is the equation that is obtained as a result of introducing into the linear second-order differential equation  $\ddot{x} + \omega_0 x = 0$  the nonlinear term  $x^3$ . It is not clear to us how the partitioning method presented will aid in telling us how the amplitude and frequency of oscillations vary with time.

Alfred A. Wolf: I am most grateful to Dr. Ku for his illuminating and encouraging discussion. In general we appear to be in agreement on all the points discussed. In particular, Dr. Ku's reference to the synthesis of nonlinear systems is of particular interest. In my dissertation,<sup>1</sup> a stability theory was given. There are several interesting consequences of this theory, one of which was reported upon previously.<sup>1,2</sup> I would now like to discuss briefly another consequence that leads to a synthesis procedure.

Suppose the Laplace transform of  $x(t)$ , the solution or response of the system, after suitable change of variable as discussed in reference 1 is obtained by the Taylor-Cauchy transform<sup>3</sup> or some other method<sup>1</sup> and given by

$$X(w) = \sum_{n=0}^{\infty} h_n w^n \tag{95}$$

where  $w$  is a complex variable corresponding to  $t$  and functionally related to the complex variable,  $s$ , of the Laplace transform and  $h_n$  denotes the coefficients of  $w^n$ .

Under certain conditions  $X(w)$  will be meromorphic. A simple test for determining this is by applying Kronecker's theorem.<sup>1,4</sup> Then  $X(w)$  is expressible as the rational fraction,

$$X(w) = \frac{N(w)}{P(w)} \tag{96}$$

where

$$N(w) = \sum_{r=0}^m a_r w^r \tag{97}$$

and

$$P(w) = \sum_{r=0}^n b_r w^r \tag{98}$$

then

$$a_r = \sum_{k=0}^r h_k b_{r-k} \tag{99}$$

such that

$$a_r \begin{cases} = 0 & \text{when } r > m \\ \neq 0 & \text{when } r \leq m \end{cases} \tag{100}$$

The upper condition of equation 100 determines the  $b_r$ 's while the lower condition

determines the  $a_r$ 's once the  $b_r$ 's are known. It is evident therefore that it is possible to calculate the zeros and poles of  $X(w)$ . Thus, the stability is determined.

# EXAMPLE

Let us determine the coefficients  $a_k$  and  $b_k$  for the following problem in terms of the coefficients  $h_n$ .

$$\sum_{n=0}^{\infty} h_n w^n = \frac{a_0 + a_1 w}{b_0 + b_1 w} \tag{101}$$

$$= \frac{a_0}{b_0} \frac{1 + \frac{a_1}{a_0} w}{1 + \frac{b_1}{b_0} w}$$

Therefore

$$a_0 = h_0 b_0 \tag{102}$$

$$a_1 = h_0 b_1 + h_1 b_0 \tag{103}$$

$$0 = h_1 b_1 + h_2 b_0 \tag{104}$$

$$0 = h_2 b_1 + h_3 b_0 \tag{105}$$

For nontrivial solution a condition on the  $h$ 's:

$$\begin{vmatrix} h_1 & h_2 \\ h_2 & h_3 \end{vmatrix} = 0 \tag{106}$$

Hence

$$\frac{b_1}{b_0} = \frac{h_2}{h_1} \tag{107}$$

$$\frac{a_0}{a_1} = h_0 \tag{108}$$

$$\frac{a_1}{a_0} = \frac{h_0 h_1}{h_0 h_2 + h_1^2} \tag{109}$$

The  $h_n$ 's are then determined in terms of the  $a_n$ 's. Turning the problem around we have the synthesis problem; that is, given the  $a_n$ 's, which specify the characteristics of a certain system, the  $h_n$ 's are calculated. The dynamics of the system would be specified by a nonlinear differential equation in which at least some of the coefficients of the differential equation are unknown. The  $h_n$ 's are then functions of these unknown coefficients. Since the  $a_n$ 's are specified, the  $h_n$ 's are determined, hence the unknown coefficients describing the system dynamics are calculated. This is one kind of synthesis procedure leading to systems with a specified degree of stability and the parameters of the system are determined in such a way as to give this desired stability.

I would like to thank Dr. Cosgriff both for his comments and for affording me the opportunity of replying to the points raised.

The two main points which Dr. Cosgriff raised deal with the method of Frobenius and its relation to the method of the paper and the generalization of expanding the auxiliary forcing function.

The method of Frobenius, which has been discussed in some detail in reference 5, is a special case of the Partition Theory. However, the Frobenius method only gives rise to power series solution, but the Partition Theory gives rise to a general class of solution forms depending on the point of parti-



As pointed out in the paper, these can be Power series, Dirichlet series, geometric series, Orthogonal series, and others. The form of solution depends on the partition point and the form of the auxiliary function, which provides flexible control on the nature of the solution's form. In the Frobenius method the solution is assumed to be a power series; in the partition scheme the auxiliary function is expanded into an analytic series, not necessarily a power series. If the auxiliary function is expanded into a power series, the solution is only a power series if the point of partition is at the highest derivative. It will also be a power series if only one term remains after partition and reordering of terms. If the partition does not allow the next order derivative to remain on the left member then care must be taken so that the resulting recurrence equation will not be awkward. In this sense, a partition at the  $n$ -th derivative, with all other linear terms transposed, is equivalent to the method of Frobenius. In the paper, the effect of changing the partition point was illustrated. This illustration showed how the solution form changed from a power series to a Dirichlet series.

We now turn to the second point of Cosgriff's discussion dealing with the connection of analytic auxiliary functions other than power series. This point is discussed in great detail in reference 1. Instead of going into such detail here an example will be given.

#### EXAMPLE

In virtue of the theory in reference 1, the following expansion of  $A(t)$  is justified.

$$A(t) = \sum_{n=0}^{\infty} c_n e^{-\lambda_n t} \quad (110)$$

In our previous examples  $A(t)$  was expanded into a power series. The recurrence relations resulting from this kind of expansion are simple and elegant. The next simplest expansion of  $A(t)$  is given by equation 110 and the simplest recurrence relations among these results when  $\lambda_n$  is chosen as

$$\lambda_n = an \quad (111)$$

Equation 1 now has the solution

$$x(t) = \sum_{n=0}^{\infty} c_n P_n(t) \quad (112)$$

where  $P_n(t)$  are the exponential moments of the folded impulsive response of the linear part given by

$$P_n(t) = \int_0^t e^{-nar} y(t-\tau) d\tau \quad (113)$$

It is evident that on expanding  $e^{-nar}$  into a power series a relation exists between the exponential moments and the moments  $Q_n(t-\tau)$ ,  $Q_n(t)$ .

By choosing  $\lambda_n$  according to equation 111, we have the special case of the Dirichlet series, known as the Power Dirichlet series. This terminology is adapted since equation 112 reduces to a power series in  $Z$  when  $Z = e^{-at}$  and  $\lambda_n$  is given by equation 111. Consider, for example, a system described by

$$\frac{dx}{dt} + fx + bx^2 = g \quad (114)$$

where  $f, b$ , and  $g$  are constants, subject to the initial conditions

$$x(0) = 0 \quad (115)$$

Equation 114 may be partitioned in two ways. The first might be at the highest derivative or a second choice might include, in addition to the derivative, the function itself. Consider the former partition first:

$$\frac{dx}{dt} = g - fx - bx^2 \quad (116)$$

Thus,

$$\frac{dx}{dt} = \sum_n c_n e^{-nat} \quad (117)$$

performing the indicated operations and equating like powers of  $Z = e^{-at}$  the following recurrence relation is obtained.

For  $n > 0$ :

$$C_n = (g - fk - bk^2)c_n^{(0)} + \frac{fc_n}{na} - bc_{n,0}^{(2)} + 2bk \frac{c_n}{na} \quad (118)$$

where  $k$  is a constant of integration.

For  $n=0$  we are enabled to determine  $k$ ; that is, letting  $n=0$  we obtain

$$bk^2 + fk - g = 0 \quad (119)$$

or

$$k = \frac{-f \pm \sqrt{f^2 + 4bg}}{2b} \quad (120)$$

For  $n=1$ :

$$c_1 = f \frac{c_1}{a} - bc_{1,0}^{(2)} + 2bk \frac{c_1}{a} \quad (121)$$

or

$$c_1 \left( \frac{f}{a} + \frac{2bk}{a} \right) = 0 \quad (122)$$

The second factor is not zero; hence,

$$c_1 = 0 \quad (123)$$

For  $n=2$  it is possible to determine the characteristic exponent,  $a$ . After simplification

$$c_2 \left[ 1 - \frac{f}{2a} + \frac{bc_1^2}{a^2} - \frac{2bk}{2a} \right] = 0 \quad (124)$$

Either or both factors may be zero. We shall assume that the second factor is zero to determine  $a$ . Noting  $c_1 = 0$  from equation 123 and substituting the value of  $k$  given by equation 120 yields the results

$$a = \pm \frac{1}{2} \sqrt{f^2 + 4bg} \quad (125)$$

This result is known to be correct since it can be checked against the solution as obtained by separating the variables. The characteristic exponent in this case has two values. Both must be used in equation 126 giving rise to Dirichlet series

$$x = \sum_{n=1}^{\infty} \frac{c_n}{-na} e^{-nat} + k \quad (126)$$

If the characteristic exponent had  $k$  values then  $k$  series would be needed in addition to the constants of integration which may arise. This procedure is similar to that used in linear differential equations. The remainder of the solution follows a similar procedure to that described elsewhere.

If the second partition scheme is used instead, a different form of  $x(t)$  is obtained; this is easily checked.

The third case proposed by Cosgriff, namely expanding

$$A(t) = \sum_{n=0}^{\infty} c_n t^n e^{-at} \quad (127)$$

evidently follows a similar procedure. Actually a more interesting case for investigation is to expand  $A(t)$  according to the relation

$$A(t) = \sum_{n=0}^{\infty} c_n t^n e^{-nat} \quad (127)$$

Regarding the question of generating transforms for these cases like the Taylor-Cauchy transform, the answer is clearly that this is indeed possible and has been done. These results may be reported upon later.

The questions raised by Dr. Kazda and Dr. Bilal are much appreciated and greatly valued. The three points raised by them are particularly pertinent to the theory of partition and I shall take up certain aspects of these points not covered before.

Let us consider the first point relative to the stability of nonlinear systems. In reference 1 it was shown that a certain class of linear systems existed (not necessarily physically realizable) such that the response,  $x_n(t)$  for  $n=0, 1, 2, \dots$ , of each of these systems formed a point set whose limit as  $n \rightarrow \infty$  is the response of the nonlinear differential equation given in the paper as equation 1. The members of the set  $x_n(t)$  are generated as follows:

$$x_0(t) = \int_0^t g(\tau) y(t-\tau) d\tau \quad (128)$$

and

$$x_n(t) = x_0(t) - \int_0^t F\{x_{n-1}(\tau)\} y(t-\tau) d\tau \quad (129)$$

for  $n \geq 1$

The limit referred to,

$$x(t) = \lim_{n \rightarrow \infty} x_n(t) \quad (130)$$

is rigorously proved in references 1 and 6 as a uniformly convergent procedure. Because of uniformity, the following theorem is proved.

*Theorem 1.* Suppose that

$$\lim_{t \rightarrow \infty} x_n(t) = 0 \quad (131)$$

then there exists an  $N > 0$  such that for each  $n \geq N$  the limit, equation 131, governs the behavior of  $x(t)$  so that

$$\lim_{t \rightarrow \infty} x(t) = 0 \quad (132)$$

What this theorem implies is that there is, from a certain member on, in the set of responses  $\{x_n(t)\}$  a linear member whose stability governs the stability of the nonlinear system. While theoretically this is an important result, the solution of  $x_n(t)$  for a given  $n \geq N$  is cumbersome and its stability is not immediately discernible. To ameliorate this situation one makes use of the following theorem.

**Theorem 2.** The members of the point set  $\{x_n(t)\}$ , under the conditions specified in the paper, are each of exponential type.

This implies that each member has a Laplace transform, and moreover, the singularities for a lumped parameter system are poles for the  $\mathcal{L}\{x_n(t)\} = X_n(s)$  for all finite  $n$ . In the limit the function  $x(t)$  is not necessarily of exponential type so it does not necessarily have a Laplace transform.

Utilizing the partition method the series solution of  $x(t)$  can be obtained. By proper partition, the power series solution, or any other series solution is obtainable. Even though the expansion of  $x(t)$  does not have a Laplace transform directly it can be obtained as near as desired by a limiting procedure on the solution of  $x_n(t)$ . Since the latter cannot easily be obtained, a solution of  $x(t)$  is found in series form. As shown in reference 1, the coefficients,  $C_n$  of this series gives the singularities of

$$\mathcal{L}\{x(t)\} = \lim_{n \rightarrow \infty} \mathcal{L}\{x_n(t)\}$$

which may not be poles. These singularities form a convex singularity hull. If this

hull is in the left half-plane of  $s$ , the Laplace complex variable, the system is stable; otherwise, it is not.

For absolute stability, the convex singularity hull must be in the left half plane for all inputs. To do this practically white noise must be applied to the system, as described in the paper, and the average position of the convex singularity hull must be noted. If this average position is in the left-plane the system is absolutely stable.

Let us now take up the second point raised by Kazda and Bilal namely, "under what conditions would partitioning at points other than the highest derivative be desirable." There are a number of instances where this is desirable. When there is interest in displaying certain harmonics of the system partitioning might be utilized so as to display these in the solution expansion as a linear combination. Another place where partitioning at a point other than the highest derivative is important is the case where the recurrence relations can be made to terminate, or at least be made simpler.<sup>1</sup> There are two important subcases where this simplification takes place. One pertains to the case of a given expansion of the auxiliary function. In this case the partition point can be varied by trial to see if the recurrence relations become simpler. The other subcase pertains to changing the form of expansion of the auxiliary function. For example, if the partition is to be made so that the highest derivative say the  $k$ -th and the  $(k-1)$ st remain in the partition, the Dirichlet series might be used as an expansion

of the auxiliary function. This would give the simplest moment function, and the solution of the system would be another Dirichlet series.

Regarding the third point of Kazda and Bilal's discussion, the auxiliary function the Dirichlet series would be selected. The point of partition would depend on how wished to display the result. If the result to show functional dependence on  $\omega_0$  in terms of harmonics, the partition would be made after the second linear term namely  $\omega_0$ . For other purposes it might be made at  $\frac{\pi}{2}$ . In either case, as shown above, the characteristic exponent would contain information relative to the frequency of oscillation. Determining the amplitude somewhat more complicated if the series does not terminate. What might be done here is the rewriting of the terms of the series as trigonometric functions then utilizing trigonometric identities to display amplitude and phase information.

#### REFERENCES

1. See reference 2 of the paper.
2. See reference 7 of the paper.
3. See reference 8 of the paper.
4. THE TAYLOR SERIES (book), P. Dienes, Oxford University Press, New York, N. Y., 1931, pp. 89-90.
5. Reply by A. A. Wolf to Bohn's note on the EQUIVALENCE OF THE TAYLOR-Cauchy LAURENT-CAUCHY TRANSFORM ANALYSIS AND CONVENTIONAL METHODS, *Proceedings, Institute of Radio Engineers*, New York, N. Y., vol. 49, no. Jan. 1961, pp. 358-61.
6. See reference 3 of the paper.

## Feedback Compensation: A Design Technique

G. J. THALER  
MEMBER AIEE

J. D. BRONZINO  
STUDENT MEMBER AIEE

D. E. KIRK  
STUDENT MEMBER AIEE

THE DESIGN OF feedback compensation for feedback control systems is a subject which is treated in the technical literature as an undertaking incidental to the design of specific systems.

Most textbooks are concerned primarily with cascade compensation; the few<sup>1-3</sup> that devote reasonable space to feedback compensation treat special cases or reduce minor loops, one at a time. None provide a systematic approach which has general applicability. None of the suggested techniques give much guidance as to suitability of a selected compensation path or compensator. Excessive labor in many applications is expended in repeated trial-and-error solutions.

The technique suggested herein is generally applicable to systems of any order, and to nonlinear systems under certain conditions. It permits design of

feedback compensators by following well-known and well-practiced techniques of cascade compensation, and offers some insight into whether selected paths and compensators are appropriate. Moreover, the labor involved is no greater than that required for cascade compensation design.

### Fundamental Principles

Numerous specific procedures apply to feedback control system design. In general, however, after the selection of what may be called "unalterable components," the system gain, or type number, or both, are fixed to satisfy specifications of static and steady-state accuracy. A stability check follows, and a need for compensation is usually indicated. When feedback compensation is to be used, a suitable path and compensator must be chosen.

Choice of a path is restricted by availability of signal pick-off points, signal feed-in (summing) points, and suitable and accurate pick-off and measuring devices. For simpler systems path selection is relatively easy; for complex systems, theoretical information is limited. Some work<sup>4</sup> has been done, but much more is needed.

The suitability of a compensator, however, is more readily judged. First must not affect the low-frequency gain of the forward path. When this requirement is applied, a feedback path, utilizing gain adjustment only, alters the forward gain at all frequencies, thus affecting the steady-state accuracy of a type-0 system in response to a step input or a load disturbance, as well as the steady-state accuracy of a type-1 system in response to a ramp input, load torque, etc. First derivative feedback affects the velocity error of a type-1 system following a ramp

Paper 61-752, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department, presented at the AIEE Summer General Meeting, Ithaca, N. Y., June 18-23, 1961. Manuscript submitted March 16, 1961; made available for printing April 18, 1961.

G. J. Thaler, J. D. Bronzino, and D. E. Kirk are with the United States Naval Postgraduate School, Monterey, Calif.



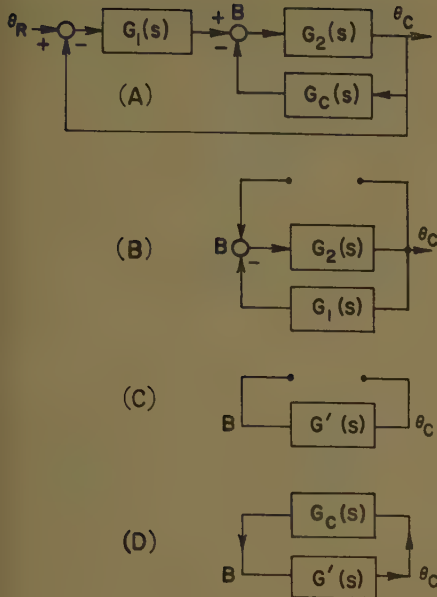


Fig. 1. Fundamental manipulation procedure

- A—Typical block diagram
- B—Break compensator path and rearrange
- C—Reduce to equivalent single block
- D—Reconnect compensator

but does not affect the steady-state response to a step or a load disturbance. Second-derivative feedback affects acceleration-lag error of a type-2 system.

The compensator must, as a second requirement, be capable of stabilizing the system and of adjusting the transient response to meet specifications. Aside from such considerations as physical realizability and economics, this simply

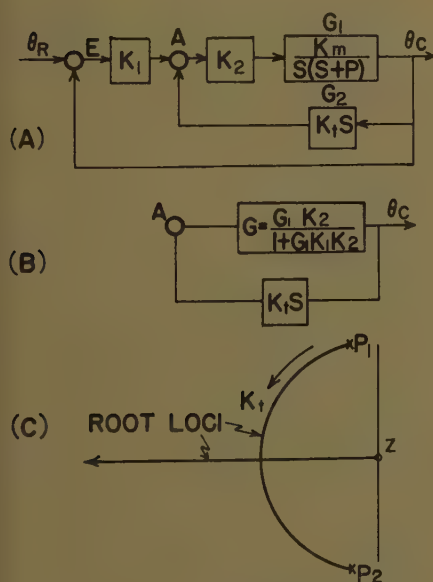


Fig. 2. Tachometer feedback

- A—Block diagram
- B—Equivalent single loop
- C—Root-locus plot

means that the compensator must be able to move the roots of the characteristic equation to suitable locations; or that it must adjust phase and gain margins and resonant frequency to acceptable values. This is equivalent to saying the modes of transient oscillation are specified by root locations only. True, zeros of the system function adjust the amplitudes of oscillation modes and thus may alter system response to the extent of making an apparently acceptable root configuration unacceptable. Conversely, a root in an undesirable location may be essentially nullified by a properly located zero. In general, selection of a compensator that will provide acceptable root locations is the reasonable starting point.

Since the technique herein developed adopts this approach, certain adjustments of the block diagram are permissible. If the characteristic equation and its roots are the sole interest, all command and disturbance inputs may be kept constant at zero for sake of convenience. This permits elimination of several summing points which otherwise would be annoying.

### Basic Block Diagram Manipulation

Assume that all gains have been set to satisfy steady-state specifications, and that a feedback compensation path has been selected for trial use. Compensation technique then depends on block-diagram manipulation. The procedure is:

1. Set all input and disturbance signals at zero, and eliminate summation points wherever possible.
2. With compensator path removed, but with its pick-off and feed-in points noted, rearrange block diagram so that feed-in and pick-off points become input and output, respectively. Reduce this diagram to a single equivalent block as in Fig. 1 (A), (B), and (C).
3. Connect compensator to pick-off and feed-in points, thus forming a single-loop block diagram as in Fig. 1(D).
4. Design compensator by using conventional cascade compensation methods.

This technique is applicable to linear systems of any complexity, with the compensator itself either linear or nonlinear. No change in conceptual approach is required if applying this method to certain classes of nonlinear systems, but there will be unavoidable increases in labor.

### Second-Order Systems

#### TACHOMETER FEEDBACK COMPENSATION

For compensating lightly damped instrument servos, tachometer feedback is

commonly used as illustrated by Fig. 2(A), where the motor-load unit is assumed to be of second order. Block-diagram manipulation produces Fig. 2(B), for which the root-locus plot is as shown in Fig. 2(C). Adjustable gain is  $K_t$ , the tachometer constant, on this plot. An increase in  $K_t$  drives the roots from the poles, increasing the damping but leaving  $\omega_n$  unchanged.

The ordinate and abscissa of this root-locus plot may be nondimensionalized to form a universal nomograph for tachometer compensation design<sup>6</sup> of static positioning systems. Instead of resorting to a nomogram, however, Appendix I shows that

$$K_t = \frac{2\zeta\omega_n - p}{K_2K_m} \quad (1)$$

where  $K_2$ ,  $K_m$ ,  $p$ , and  $\omega_n$  are defined by the uncompensated system, and  $\zeta$  is specified by the desired root location.

Two ways of using this equation are apparent. Given the values for motor pole and gain, and with the threshold error specified, the required loop gain  $K_1K_2K_m$  is determined. If the gain division into  $K_1$  and  $K_2K_m$  is made arbitrarily, then  $K_t$  is evaluated from the equation and tachometer circuit designed.

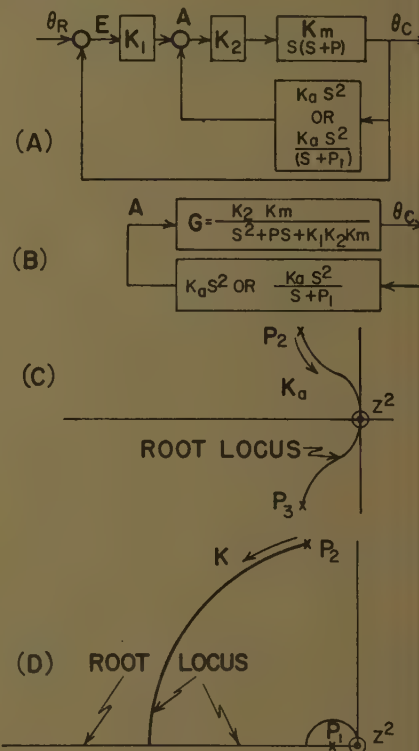


Fig. 3. Acceleration feedback, exact and approximate

- A—Block diagram
- B—Equivalent single loop
- C—Root locus for  $K_a s^2$ , or for  $K_a s^2/(s+p)$  if  $p_1 \gg p$
- D—Root locus for  $K_a s^2/(s+p_1)$  if  $p_1 < p$

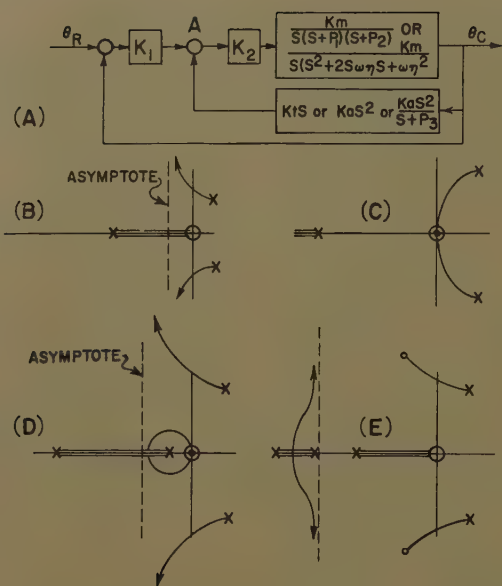


Fig. 4 (left). Third-order systems

- A—Block diagram  
B—Root locus for tachometer feedback  
C—Root locus for acceleration feedback  
D—Root locus for approximate acceleration feedback  
E—Root locus for feedback through  $\frac{K_a s^2 (s^2 + \alpha s + \beta)}{(s + p_1)(s + p_2)}$

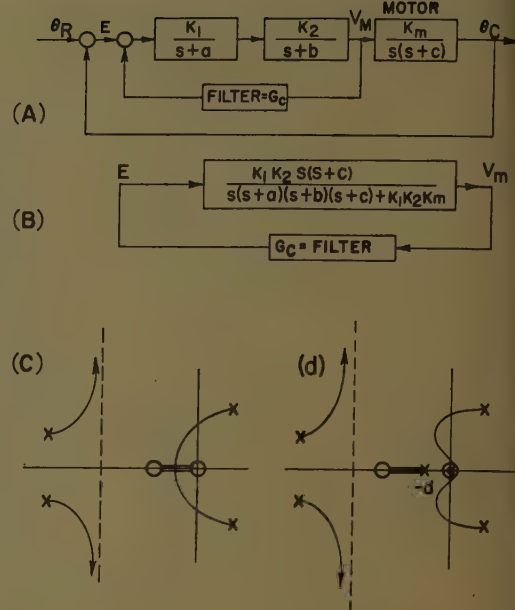


Fig. 5 (right). A fourth-order system

- A—Block diagram  
B—Equivalent single loop  
C—Root locus for  $G_c = K$   
D—Root locus for  $G_c = K_a s / (s + d)$

On the other hand, if a tachometer of known  $K_t$  is used, the equation permits evaluation of  $K_2 K_m$  and thus specifies the gain subdivision which permits use of full  $K_t$ , or any selected fraction thereof. For example, if  $K_m = 200$ ,  $p = 23$ , and threshold accuracy requires that  $K_1 K_2 K_m = 10^5$ , then it may be decided arbitrarily that  $K_1 = 100$  and  $K_2 K_m = 10^3$ . In this case, for a compensated system with  $\zeta = 0.7$ , the gain must be

$$K_t = \frac{2\zeta\omega_n - p}{K_2 K_m} = \frac{2(0.7)10^{5/2} - 23}{10^3} = 0.42$$

On the other hand, if  $K_t$  is 0.32 volt per radian per second, and this full value is used, then the necessary gain subdivision in the forward path is found by noting that

$$K_2 K_m = \frac{2\zeta\omega_n - p}{K_t} = \frac{2(0.7)10^{5/2} - 23}{0.32} = 1,310$$

$$K_1 = \frac{K_1 K_2 K_m}{K_2 K_m} = \frac{10^5}{1,310} = 76.3$$

$$K_2 = \frac{1,310}{200} = 6.5$$

#### ACCELERATION FEEDBACK COMPENSATION

Since tachometer feedback increases lag error in following a ramp, acceleration feedback is often used to damp transient oscillations of second-order systems. The basic block diagram, Fig. 3(A), after manipulation, takes the form of 3(B). Two transfer functions are shown for the compensator,

$K_a s^2$   
for pure acceleration feedback, and  
 $K_a s^2 / (s + p_1)$

for the commonly used tachometer followed by a lead network.

Note that approximate acceleration feedback (using the filter) actually provides much better control of dynamic performance than pure acceleration feedback, if the pole is suitably located. Fig. 3(C) shows the root locus for pure acceleration feedback. Damping ratio  $\zeta$  cannot be improved, although decreased bandwidth is possible. If the tachometer-filter combination has the pole far out on the negative real axis, little improvement is available. But, with the pole near the origin, the root locus, as seen in Fig. 3(D), permits improved damping with large or small  $\omega_n$ , depending on the amount of feedback gain. Although the latter placement of the pole introduces a real root near the origin for small gain, the compensator pole is a zero of the system function. For a compensator pole less than 0.15 times the motor pole, this zero effectively cancels the residue at the real pole if the feedback gain is small enough to keep the complex roots on the outer semicircle. If the gain is large enough to place complex roots on the inner circle, the real root is large enough to neglect. In either case, second-order response obtains; there is no tail to the step response due to the presence of the real root.

For the case of acceleration feedback using the function

$$K_a s^2 / (s + p_1)$$

another simple design nomograph<sup>6</sup> can be constructed which is applicable to most cases of interest. For acceleration feedback, also, an equation is available for computation of feedback gain (see Appendix II):

$$K_a = \frac{1}{K_2 K_m} \left[ a\omega_{nf} - p - p_1 + \frac{aK_1 K_2 K_m}{\omega_{nf}} + \frac{ap_1}{\omega_{nf}} - \frac{p_1 K_1 K_2 K_m (a^2 - b^2)}{\omega_{nf}^2} \right] \quad (2)$$

where  $\omega_{nf}$  is the natural frequency of the final or compensated roots,

$$a\omega_{nf} = \zeta f\omega_{nf}$$

is the real part of the root value and the other symbols are defined in Fig. 3. This equation reduces to a simpler form if the desired root is to be located on the outer circular arc of Fig. 3(D)

$$K_a \approx \frac{1}{K_2 K_m} \times \left[ 2a\omega_{nf} - p - p_1(1 + a^2 - b^2) + \frac{ap_1}{\omega_{nf}} \right] \quad (3)$$

and, as an approximation, the last two terms may be dropped, since they are small when  $\omega_{nf}$  is large, giving

$$K_a \approx \frac{2a\omega_{nf} - p}{K_2 K_m}$$

These equations are readily applied. Assume that the system previously considered is to be compensated with approximate acceleration feedback instead of tachometer feedback. Then, for roots on the outer semicircle

$$\omega_{nf} = \sqrt{K_1 K_2 K_m} = \sqrt{10^5} = 316; K_2 K_m = 10^3; p = 23;$$

choose  $p_1 = 2.3$ ;  $\zeta = 0.7$ ; therefore

$$a = 216 / \sqrt{2} = 223.$$

$$K_a = \frac{2a\omega_{nf} - p}{K_2 K_m} = \frac{446 - 23}{10^3} = 0.423$$

On the other hand, if the system band-



width is to be reduced, the exact equation must be used, and these values apply

$$K_1 K_2 K_m = 10^5; K_2 K_m = 10^3; p = 23;$$

choose  $p_1 = 6$ ; then for  $\zeta = 0.7$ ,  $a = b = 6$ , and  $\omega_{nf} = 6\sqrt{2}$ . It follows that

$$K_a = \frac{1}{10^3} \times \left[ 36\sqrt{2} - 23 - 6 + \frac{6 \times 10^5}{6 \times 2} + \frac{(6)(23)(6)}{6 \times 2} - 0 \right] \\ = 10^{-3} [36\sqrt{2} - 29 + 0.7 \times 10^5 + 0.7(138)] \\ \cong 700 \quad (4)$$

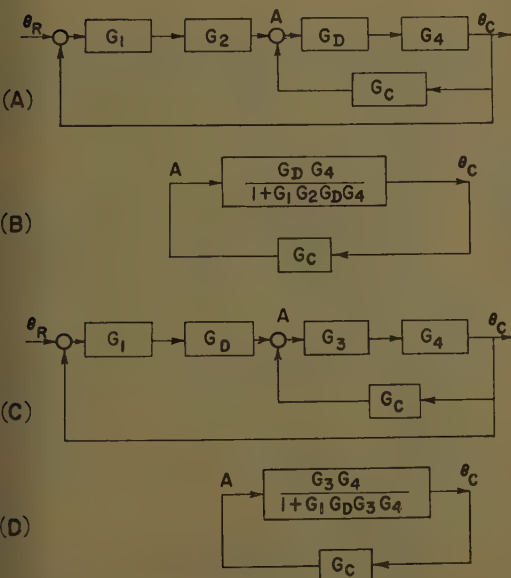
Note from this illustration that a good approximation exists when  $K_1 K_2 K_m$  is large and  $\omega_{nf}$  is on the inner circle

$$K_a \cong a K_1 / \omega_{nf} \quad (5)$$

### Third-Order Systems

Most control systems are at least of the third order, even if the real root may be large enough to neglect, in which case approximation by second-order-system calculations may be adequate. When the real root is not large, analysis of feedback compensator effects is greatly aided by the technique under discussion. Fig. 4(A) shows the block diagram for a case where the loop-transfer function of the minor loop contains all poles of the forward-transfer function; i.e., the compensation signal is fed back around all poles of the forward path.

Forward-path poles may be either real or complex. For the high gains usually required, complex roots are normally in the right-half plane. Figs. 4(B), (C), and (D) are root-loci sketches for different feedback compensators. Clearly, pure acceleration feedback will not work; tachometer feedback may work if the



asymptote is in the left-half plane, but then the real root will probably dominate. Approximate acceleration feedback may also work if the asymptote is in the left-half plane, in which case a pair of complex roots with relatively low  $\omega_n$  probably can be made to dominate.

None of these compensation schemes is satisfactory for a variety of applications. Much time is saved, therefore, if the manipulation technique is used in a preliminary examination of all suggested compensators to eliminate unsuitable schemes. Fig. 4(E) shows a scheme where the feedback function generates suitably located complex zeros, allowing better placement of the complex roots. Whether the roots can be made dominant or not depends on numerical values involved and on the specific system function obtained.

### Higher-Order Systems

This manipulation technique may be applied to linear systems of any order. Choice of compensation path depends on application, but the path need not enclose all forward-path poles. This leads to pole-zero configurations with numerous zeros. The choice of compensator, while restricted by steady-state specifications, may be aided by inspecting the uncompensated root locus.

Fig. 5(A) is the block diagram of a fourth-order positioning system, for which the proposed compensation is attempted by feeding back the motor-armature voltage through a filter. Fig. 5(B) is the

manipulated block diagram and Fig. 5(C) shows the root locus for the feedback loop connected with  $G_c = K$  but without filter. The connection obviously would stabilize the system, but would disturb the steady-state velocity-lag error. Fig. 5(D) shows the addition of a filter with transfer function  $K_s s / (s + d)$ , which provides a possible compensator, but may not permit satisfactory root location. The next step is to try a more complicated filter, perhaps with complex zeros. Or, a new compensation path may be indicated.

### Nonlinear Systems

If a system is nonlinear and linear feedback compensation is to be used, analysis and design may be accomplished using the describing function representation of the nonlinear component. Block diagram, Fig. 6(A), shows a system wherein the compensator  $G_c$  feeds around the nonlinearity  $G_D$ . After block-diagram manipulation, the equivalent single loop is as shown in Fig. 6(B), and algebraic manipulation provides the stability relationship

$$G_4(G_c + G_1 G_2) = -1/G_D \quad (6)$$

In like manner, Fig. 6(C) shows a system with nonlinear element outside the compensating loop; Fig. 6(D) is the equivalent single-loop diagram, and again the stability relationship is produced by algebraic manipulation

$$-\frac{1}{G_D} = \frac{G_1 G_2 G_4}{1 + G_1 G_2 G_4} \quad (7)$$

Fig. 6 (left). Linear compensation of a nonlinear system

- A—Block diagram, compensation feedback around nonlinearity
- B—Equivalent single loop
- C—Block diagram, compensation feedback around linear elements only
- D—Equivalent single loop

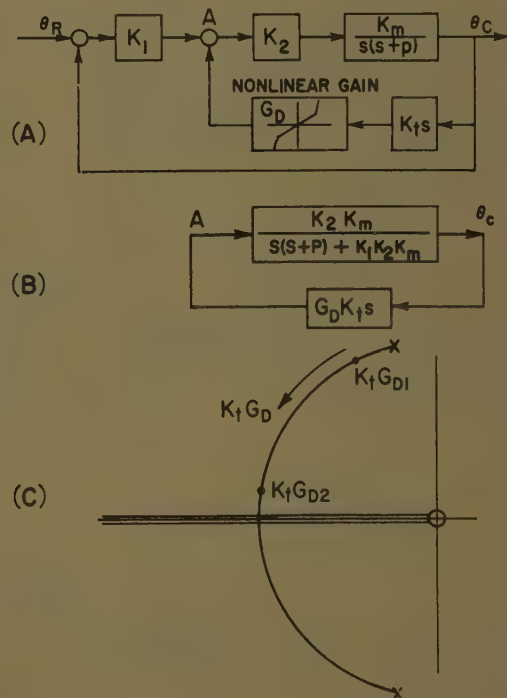


Fig. 7 (right). Nonlinear compensation of a linear system

- A—Block diagram
- B—Equivalent single loop
- C—Root-locus diagram

Equation 6 and equation 7 may be analyzed for stability by Nyquist methods or by the root-locus method if the describing function is amplitude-sensitive only.

If nonlinear compensator  $G_{DC}$  is used with a linear system, then this technique is particularly helpful. Fig. 7(A) is the block diagram of a second-order servo, which is lightly damped to provide fast rise time and small velocity-lag error. If heavy damping is needed for large disturbances, tachometer feedback through a nonlinear amplifier can be used. The nonlinearity is adjusted to give proportionate feedback gain—high for large disturbances; low or zero for small. Using the equivalent single-loop diagram in Fig. 7(B), and the root locus in Fig. 7(C), it is easy to design the nonlinearity with small-disturbance roots at  $K_t G_{D1}$ , and with large disturbance roots at  $K_t G_{D2}$ .

## Conclusions

Design of feedback compensation to control root location, which, in turn, controls stability and dynamic performance, can be reduced to an equivalent problem in cascade compensation by a simple manipulation of the block diagram. In conjunction with this, the root-locus method provides a particularly convenient means of analysis and design. These methods permit derivation of such simple equations as for the tachometer, and also simplify the logical choice of compensation paths and compensators for complex systems.

As some nonlinear systems also may be treated with this technique, nonlinear compensator design thus may be made easier in such cases. Linear compensator design for nonlinear systems is feasible, too, and labor is reduced in the process when considered in comparison with other methods.

## Appendix I

### Derivation of Tachometer Feedback Relationships

From Fig. 2(B)

$$K_t s \left( \frac{K_2 K_m}{s(s+p) + K_1 K_2 K_m} \right) = -1$$

$$K_t = - \left( \frac{s(s+p) + K_1 K_2 K_m}{K_2 K_m s} \right) \quad (8)$$

Let the desired root be at  $s = (-a + jb)\omega_n$ , and note that  $a^2 + b^2 = 1$ . Substitute for  $s$ , and consider only the real part of the result, since  $K_t$  is a real number

$$K_t = - \left\{ \frac{[\omega_n^2(1+a^2-b^2) - a p \omega_n](-a) + b^2(p\omega_n - 2a\omega_n^2)}{K_2 K_m \omega_n (a^2 + b^2)} \right\}$$

$$= \frac{2a^3\omega_n^2 + 2ab^2\omega_n^2 - a^2 p \omega_n - b^2 p \omega_n}{K_2 K_m \omega_n}$$

$$= \frac{2a\omega_n - p}{K_2 K_m} = \frac{2\xi\omega_n - p}{K_2 K_m} \quad (1)$$

where  $a = \xi$  since  $-a\omega_n = -\xi\omega_n$  is the real part of the root number.

An alternate and even simpler derivation obtains from the characteristic equation:

$$s^2 + (K_2 K_m K_t + p)s + K_1 K_2 K_m = 0 \quad (9)$$

and

$$K_2 K_m K_t + p \triangleq 2\xi\omega_n \quad (10)$$

thus

$$K_t = \frac{2\xi\omega_n - p}{K_2 K_m} \quad (1)$$

## Appendix II

### Derivation of Acceleration Feedback Relationships

From Fig. 3(B)

$$\frac{K_a K_2 K_m s^2}{(s+p_1)(s^2 + ps + K_1 K_2 K_m)} = -1$$

$$K_a = - \frac{(s+p_1)(s^2 + ps + K_1 K_2 K_m)}{K_2 K_m s^2} \quad (11)$$

Let  $s = (-a + jb)\omega_n$ , where  $\omega_n$  is the final or selected natural frequency for the complex roots. Substitute, expand, and consider the real part only

$$K_a = \frac{1}{K_2 K_m} \left( a\omega_n - p - p_1 + \frac{(K_1 K_2 K_m + p p_1)a}{(a^2 + b^2)\omega_n} - \frac{K_1 K_2 K_m p_1 (a^2 - b^2)}{(a^2 + b^2)\omega_n^2} \right) \quad (12)$$

but  $a^2 + b^2 = 1$ , thus

$$K_a = \frac{1}{K_2 K_m} \left( a\omega_n - p - p_1 + \frac{a K_1 K_2 K_m}{\omega_n} + \frac{a p p_1}{\omega_n} - \frac{p_1 K_1 K_2 K_m (a^2 - b^2)}{\omega_n^2} \right) \quad (2)$$

$K_1$ ,  $K_2$ ,  $K_m$ ,  $p$ , and  $p_1$  are known from the system components, while  $a$ ,  $b$ , and  $\omega_n$  are specified by selecting the desired root location. If it is on the inner circle of the root locus, then the foregoing equation must be used; if on the outer circular arc,  $K_1 K_2 K_m = \omega_n^2$  so that

$$K_a \cong \frac{1}{K_2 K_m} \left( a\omega_n - p - p_1 + a\omega_n + \frac{a p p_1}{\omega_n} - p_1 (a^2 - b^2) \right) \quad (3)$$

which indicates that  $p_1$ ,  $p_1(a^2 - b^2)$ , and  $a p p_1 / \omega_n$  are small and may be neglected. Then

$$K_a \cong \frac{1}{K_2 K_m} (2a\omega_n - p)$$

An alternate derivation of the latter equation is

$$K_a K_2 K_m + p + p_1 = \sum \text{roots} = 2a\omega_n + r_3 \quad (13)$$

from which is taken

$$K_a = \frac{2a\omega_n - p + r_3 - p_1}{K_2 K_m} \quad (14)$$

For  $r_3 \cong p_1$ , with both  $p_1$  and  $r_3$  small and  $\omega_n$  large, this readily reduces to the foregoing unnumbered equation. By combining equations 14 and 2, an exact expression for the real root location is given.

$$r_3 = -a\omega_n + \frac{a K_1 K_2 K_m}{\omega_n} + \frac{a p p_1}{\omega_n} - \frac{p_1 K_1 K_2 K_m (a^2 - b^2)}{\omega_n^2} \quad (15)$$

Combining equations 14 and 3 gives an approximate value

$$r_3 \cong p_1 \left( \frac{a p}{\omega_n} - a^2 + b^2 \right) \quad (16)$$

## References

1. INTRODUCTION TO THE DESIGN OF SERVO-MECHANISMS (book), J. L. Bower, P. M. Schultheiss. John Wiley & Sons, Inc., New York, N. Y., 1958.
2. CONTROL SYSTEM ANALYSIS AND SYNTHESIS (book), J. J. D'Azzo, C. H. Houpis. McGraw-Hill Book Company, Inc., New York, N. Y., 1960.
3. ELEMENTS OF SERVOMECHANISM THEORY (book), G. J. Thaler. McGraw-Hill Book Company, Inc., 1955.
4. AN ANALYTICAL TECHNIQUE FOR THE DESIGN OF MULTI-LOOP SERVOMECHANISM COMPENSATION, J. D. Beecher, A. M. Pride. M. S. Thesis, Massachusetts Institute of Technology, Cambridge, Mass., May 1960.
5. FIRST DERIVATIVE FEEDBACK COMPENSATION OF CONTROL SYSTEMS, D. E. Kirk. M. S. Thesis, United States Naval Postgraduate School, Monterey, Calif., June 1961.
6. COMPENSATION OF CONTROL SYSTEMS UTILIZING ACCELERATION FEEDBACK, J. D. Bronzino. M. S. Thesis, United States Naval Postgraduate School, June 1961.

## Discussion

Thomas J. Higgins (University of Wisconsin, Madison, Wis.): This paper gives an easily grasped account—from the root-locus approach—of procedure for analyzing control systems whose steady state and transient performance may be improved by inserting a compensating network in the feedback rather than the forward link. No textbooks on basic control theory now in print contain similar content, although such information should be enfolded in every introductory study course utilizing such a text.

The authors have filled a gap in both educational and practical aspects of basic control engineering theory. Their clear exposition, well-buttressed by numerical examples, illustrates the ways in which basic theoretically derived relationships are applied.

J. J. D'Azzo and C. H. Houpis (Air Force Institute of Technology, Wright-Patterson



Air Force Base, Ohio): The technique presented in this paper is a definite contribution in the area of feedback compensation.

Some comments in the opening paragraph, however, regarding reference 2, are incorrect. Sections 14-8 and 14-9 of that textbook explain a feedback compensation technique (Method 1) with a systematic approach and limited general usefulness. Some insight also is provided into feedback compensator suitability.

The method presented in the text has two advantages over that given in the paper:

1. Output and input points are maintained.

2. The selected compensation's effect on gain, and therefore on steady-state accuracy, can be determined with greater ease. Both methods are handicapped in that they may require the solution of a third- or higher-order polynomial. Using the technique described in the paper, polynomial complexity depends on the original system's complexity. Combining this system into one transfer function with factored numerator and denominator can be tedious. Polynomial complexity depends on feedback compensator complexity when using the referenced method, and this, likewise is undesirable.

**W. C. Schultz** (Cornell Aeronautical Laboratory, Inc., Buffalo, N. Y.): The authors apparently have a good scheme for quickly determining the preliminary design of a feedback compensator network. This conclusion is circumscribed by the word "apparent," since the procedure is not clearly outlined. Hence, this question: How does the new technique differ from that proposed earlier by Ross, Warren, and Thaler?<sup>1</sup> The value of both papers might be increased if some statements, and possibly examples, were given on how to use both methods in combination.

A second question concerns implied statements that this method is generally applicable when nonlinearities are encountered. A bit more caution would seem necessary if dealing with nonlinear rather than linear systems. Would the authors please comment on this point?

#### REFERENCE

1. DESIGN OF SERVO COMPENSATION BASED ON THE ROOT LOCUS APPROACH, E. R. ROSS, T. C. WARREN, G. J. THALER. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 79, Sept. 1960, pp. 272-77.

**G. J. Thaler, J. D. Bronzino, D. E. Kirk:** We thank Professors Higgins, D'Azzo, and Houpis, and Dr. Schultz for their interest and discussions, and the latter for his additional comments in conversation

and correspondence, which led to some significant modifications in our wording.

Surprisingly, the concepts in this paper have not been exploited heretofore, even though, as Professor Higgins points out, the material is basic and important. Certainly we shall incorporate it in introductory control courses.

With regard to comments by Professors D'Azzo and Houpis and Dr. Schultz, we feel that somehow an erroneous impression has been conveyed; that is, that our research has uncovered a design procedure which can be specified, step by step. This is not correct, and we are sorry if any misunderstanding has arisen. We deliberately used the word "Technique" in our title. The paper itself is primarily descriptive of block diagram manipulation as a short-cut way of reducing feedback compensation problems to cascade compensation problems, thus permitting application of any well-developed cascade compensation procedure.

The technique does not select the best feedback path or type of compensator (except by trial and error), nor does it propose any specific form of numerical or graphical calculation. In our experience, the trial-and-error analysis (leading to selection of best paths and best compensators) has been most convenient and illuminating when the root-locus method is used. This is a personal preference, however. Bode diagrams will do the same job and may be favored by many designers. With either tool, the first step is to analyze the uncompensated system and determine whether or not compensation is needed, and if so, to what extent.

Results of such analyses are not discarded, but are put to further use in conjunction with our technique. The "equivalent single block," to which the uncompensated system is reduced by block-diagram manipulation, is always the closed-loop system function, multiplied by some quantity (constant or transfer function). Thus, if a root-locus analysis has been used originally, the uncompensated system's roots are the poles of the equivalent block, and the zeros (if any) must be calculated. The compensator block merely represents additional poles and zeros which are arbitrarily placed at selected locations to force the compensated system's roots to lie where we want them. Thus, sketches may be used, and the compensator poles and zeros moved around until the root locus looks satisfactory. This merely establishes the TYPE of compensator; thereafter, numerical design proceeds as usual with standard tools. If the original analysis used Bode diagrams, the "equivalent single block" represents closed-loop frequency response, multiplied by some function, and the resultant Bode plot is easily computed. The compensator block is just a cascaded

transfer function, and all Bode diagram design techniques apply.

As to specific comments by Professors D'Azzo and Houpis, the illustration in their textbook deals with pure tachometer feedback around a third-order motor, and there is an amplifier (gain only) between the error detector and the point at which the tachometer signal is added. Their method is to formulate the characteristic equation and rearrange it into the form  $1+F(s)=0$  where  $F(s)$  is chosen so that a parameter ( $K_t$  in this case) is conveniently located in the numerator. The original input and output terminals are retained because the denominator of  $F(s)$  is the same as that of the uncompensated loop  $G(s)$ . In comparison with the tachometer feedback illustration in the paper, the D'Azzo-Houpis manipulation has some advantages. However, our tachometer feedback illustration was merely an example of applying the technique. If the amplifier in the D'Azzo-Houpis system contains a filter, or if the compensator is not pure tachometer feedback, but also contains a cascaded filter, then their approach does not necessarily retain the original input-output points, nor does it isolate the tachometer  $K_t$  conveniently. Certainly it is not of assistance in choosing values for the pole and zero in the filter compensator. On the other hand, the technique developed in the paper retains the compensator block as a cascade unit and offers insight into the problem of compensator design.

With regard to Dr. Schultz's first question, this is easily answered in view of the preceding discussion. The Ross-Warren-Thaler procedure complements the technique developed in this paper. After selecting compensator path and type by the technique explained in this paper, the Ross-Warren-Thaler method may take over the task of carrying out the numerical design.

Dr. Schultz's second question is more difficult to answer. First, root-locus methods probably will not be as useful in nonlinear as in linear cases, despite the fact that they do apply nicely to the illustration in the paper. Frequency-response methods are likely to be more effective in general. Second, if the nonlinearity is in the compensation path, or if the compensator itself is nonlinear, the technique applies in the usual manner and we are aware of no restrictions or difficulties except those usually encountered in describing function work. In general, when the nonlinear element is in the uncompensated system and in some path other than the compensator path, then the proposed technique merely serves as a means of deriving the transfer function equations and not much more. In some cases the result may assist in the compensator design, and in others it may not. Many additional explorations could be made in this area.

# Can Electric Actuators Meet Missile Requirements?

G. C. NEWTON, JR.  
MEMBER AIEE

R. W. RASCHE  
NONMEMBER AIEE

ONE OF THE vital components of a missile is the device used to position the thrust directors or control surfaces. In addition to satisfying rather stringent dynamic performance requirements, this device must be highly reliable. Missile designers favor hydraulically operated actuators because of their high dynamic performance capabilities, relatively low weight, and small size.<sup>1</sup> However, the hydraulic actuators have still not reached the desired reliability because of sensitivity to dirt particles and other factors.

A number of other approaches have been considered for missile control surface actuation in addition to hydraulic actuators. Clutch-type servomechanisms have been developed but subsequently abandoned because of unsatisfactory reliability. A number of gas-operated fluid power systems have been considered and some of these are under development.<sup>1</sup> Among these is a gas-operated reaction-jet servomotor.<sup>2</sup> Potentially this device may offer better reliability, although its dynamic performance capabilities do not appear to be too promising.

Another alternative to hydraulic actuators is the electric motor. With the advent of solid-state control devices for electric energy, it should be possible to develop a highly reliable servo system for control surface actuation based upon an electric motor drive.<sup>3</sup> However, up to the present time, electric motors of acceptable sizes and weights have been unable to meet the dynamic performance requirements of these applications.

This paper examines the theoretical

feasibility of developing unconventional electric devices that can meet the imposed performance requirements of missile control surface applications. First, the performance parameter that is most definitive in determining the fitness of an electric motor in a high-performance application is identified. This parameter is the torque-squared-to-inertia ratio or power rate of the motor. Next, the power rate requirements of typical missile applications are established and it is shown that conventional motors of reasonable size and weight are unsatisfactory. This leads to an examination of the theoretical limit on the power rate of moving-conductor d-c motors. It is found that it is possible, theoretically, to build an electric actuator that will meet missile requirements. Fortified by this result, a tentative design of a servomotor is set forth. The paper closes with a brief examination of the prospects for a reluctance-type motor in missile applications.

## Power Rate as a Motor Performance Parameter

Fig. 1 shows schematically a servomotor connected to its load through gearing (or linkage). In order to answer the following questions:

1. What is the proper gear ratio  $R$  that should be used between the motor and load in order to minimize the size of motor required?
2. What is the minimum motor size that can be used?

certain assumptions must be made concerning the motor, load, and gearing.

The motor's controlling system is assumed to be capable of compensating for dynamic lags in the buildup of motor torque so that dynamic lags do not have to be considered in the determination of the motor size. The motor is assumed to have limits on the maximum torque,  $T_M$ , and on the maximum angular velocity,  $v_M$ , that it can produce. Furthermore, it is assumed that the peak torque and maximum velocity are not interrelated. That is, the motor is assumed to be capable of producing the peak torque at any velocity including maximum as well as zero velocity. This is not an unreasonable assumption

with suitable control, since the torque limit is related to magnetic saturation. In addition to the velocity and torque limits the motor is assumed to be characterized by a moment of inertia,  $J_M$ .

Several simplifying assumptions are made with respect to the load and gearing. The load is characterized by a peak torque requirement,  $T_L$ , a peak acceleration requirement,  $a_L$ , a peak velocity requirement,  $v_L$ , and a moment of inertia,  $J_L$ . It is assumed that the required load accelerations are independent of the load torque and velocity. Thus, the peak acceleration can be demanded at velocities very close to the peak velocity. It is assumed that the gearing has little kinetic energy in relation to the combined kinetic energies of the motor and load; thus, the inertia of the gearing may be ignored or accounted for by slight modifications of the motor and load inertias. In other respects, such as compliance and backlash, the gearing is assumed to be ideal, which means that it can be characterized by a single parameter, namely, the gear ratio,  $R_G$ , which is larger than one when the motor velocity is larger than the load velocity.

Starting with these assumptions it is seen that there are two basic relationships which must be fulfilled if the motor is to be capable of driving the load: the motor velocity capability must be equal to or greater than that required by the load, thus,

$$v_M \geq R_G v_L \quad (1)$$

and the motor torque capability must exceed the peak torque required by the load after allowance has been made for the torque required to accelerate the motor itself. In terms of the previously defined symbols this may be written as

$$T_M \geq J_M R_G a_L + \frac{J_L a_L + T_L}{R_G} \quad (2)$$

On the basis of these two relationships, Newton<sup>4</sup> discussed how the required motor size for driving a specified load can always be determined where either one or both of the conditions are effective constraints in establishing the motor size. From the results of this paper it can be shown that the second relation governs motor size when the ratio of the peak

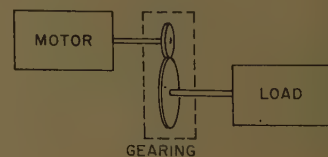


Fig. 1. Motor coupled to load through gearing

Paper 61-711, recommended by the AIEE Feedback Control Systems Committee and approved by the AIEE Technical Operations Department for presentation at the AIEE-AIEE-IRE-ISA Joint Automatic Control Conference, Boulder, Colo., June 28-30, 1961. Manuscript submitted January 13, 1961; made available for printing May 5, 1961.

G. C. NEWTON, JR., and R. W. RASCHE are both with the Massachusetts Institute of Technology, Cambridge, Mass.

The research of this paper has been supported in part by the Office of Naval Research through support extended to the Electronic Systems Laboratory, Electrical Engineering Department, Massachusetts Institute of Technology under Contract NOnw-1841(53). The authors would like to thank the members of the research department of the United Shoe Machinery Corporation who were consulted on a high-performance servomotor problem.



Table I. Representative Actuator Requirements for Flap-Controlled Missiles

	Requirement A	Requirement B
Inertia of flap assembly.....	$0.68 \times 10^{-3} \text{ kg-m}^2*$	$1.09 \times 10^{-3} \text{ kg-m}^2*$
Maximum hinge moment.....	102 newton-meters.	51.5 newton-meters
Maximum hinge acceleration.....	1,100 rad/sec <sup>-2</sup> †	625 rad/sec <sup>-2</sup> †
Maximum hinge rate.....	20 rad/sec <sup>-1</sup> †	10 rad/sec <sup>-1</sup> †

\* kilogram-meter.

† radians per second.

motor velocity to the peak unloaded motor acceleration equals or exceeds by one half a corresponding ratio for the load. That is, if

$$\frac{v_M J_M}{T_M} \geq \frac{1}{2} \frac{v_L}{a_L} \quad (3)$$

then relation 2 determines the motor size required to drive the load. For electric motors driving missile control surfaces the load acceleration requirements usually are such that condition 3 is fulfilled. In the remainder of this paper this is assumed to be so.

With condition 3 satisfied so that relation 2 governs the motor size, it is known that the motor characteristic which measures its capabilities for driving the load is its torque-squared-to-inertia ratio. This can be seen from relation 2 by finding the gear ratio,  $R$ , which minimizes the motor torque required. This optimum gear ratio, found by minimizing the right member of relation 2 with respect to  $R$ , is given by

$$R_{\text{opt}} = \sqrt{\frac{J_L a_L + T_L}{J_M a_L}} \quad (4)$$

This gear ratio may be thought of the one that matches the motor to the load impedance. Substituting this optimum gear ratio into relation 2 shows that the peak torque the motor must have, if the motor is to be capable of driving the load, is given by

$$T_M \geq 2\sqrt{(J_L a_L + T_L) J_M a_L} \quad (5)$$

The minimum permissible peak motor torque results when the equal sign in this relation holds.

In order to obtain a relation without motor parameters on the right side, condition 5 is written as

$$\dot{P}_M \geq 4(J_L a_L + T_L) a_L \quad (6)$$

where

$$\dot{P}_M \triangleq \frac{T_M^2}{J_M} \quad (7)$$

The torque-squared-to-inertia ratio of the motor is the product of the peak motor torque and the peak motor acceleration under no-load conditions. Thus, the torque-squared-to-inertia ratio is a power rate and is therefore indicated by  $\dot{P}_M$ . Relation 6 shows that the peak power rate of the motor must be four times the peak power rate of the load, or greater, for the motor to be capable of driving the load. Harris,<sup>5</sup> among others, has pointed out the significance of the torque-squared-to-inertia ratio or power rate as a servo-motor parameter.

To answer question 2, the peak power rate, given by the right side of relation 6, is computed for the load and multiplied by a factor of 4, which establishes the least value of peak power rate the motor may have in order to drive the load. From tabulated data for lines of motors, the motor which meets the required power rate can be selected. Then the gear ratio can be determined by relation 2 by replacing the inequality sign with an equal sign, and an upper and lower limit on the gear ratio, which may be used in coupling the selected motor to the load, is found. There will be a wider range of possible gear ratios as the margin by which the power rate of the motor exceeds the minimum required power rate increases.

## Actuator Requirements and Performance Levels of Available Motors

The limitations of conventional motors used as missile actuator can be demonstrated by considering the two sets of actuator performance requirements, for high-performance flap-controlled missiles, listed in Table I.

To determine the required actuator power rate, the actuator load may be considered to be an inertia and a load torque (hinge moment). The hinge moment of a control surface for a given vehicle velocity generally is an increasing function of the angular deflection, and the maximum moment tends to occur at maximum angle. Because maximum acceleration may be demanded in the vicinity of maximum deflection and maximum angular rate, the hinge moments quoted in Table I are the maximum anticipated values. In both cases shown in Table I, the maximum hinge rate is sufficiently low, in relation to the required acceleration, so that the assumption that there is no actuator speed saturation is justified; that is, relation 3 holds for most, if not all, electric motors.

Using equation 6, the minimum required motor power rates for the two sets of requirements are found to be

$$\dot{P}_{MA} = 4[(0.68 \times 10^{-3})1,100 + 102](1,100) \times 10^{-3} = 452 \text{ kw/sec} \quad (8)$$

$$\dot{P}_{MB} = 4[(1.09 \times 10^{-3})(625) + 51.5](625) \times 10^{-3} = 130.4 \text{ kw/sec} \quad (9)$$

A number of conventional, commercially available electric actuators have been analyzed with respect to power rate at the high overload condition corresponding to a 4% duty cycle. (The power rate increases with decreasing duty cycle.) The results of this analysis are summarized in Figs. 2(A) and 2(B); these are plots of power rate as a function of continuous horsepower rating and motor weight. Requirement A is met by two of the considered actuators, namely, the

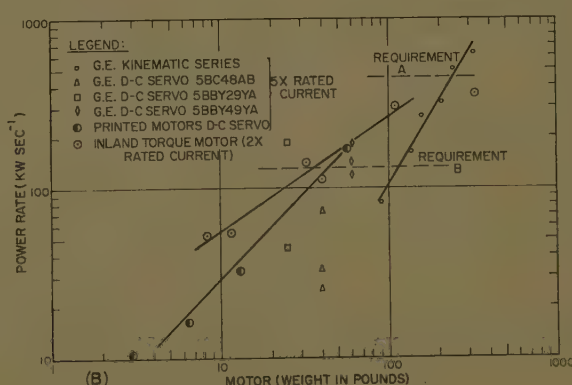
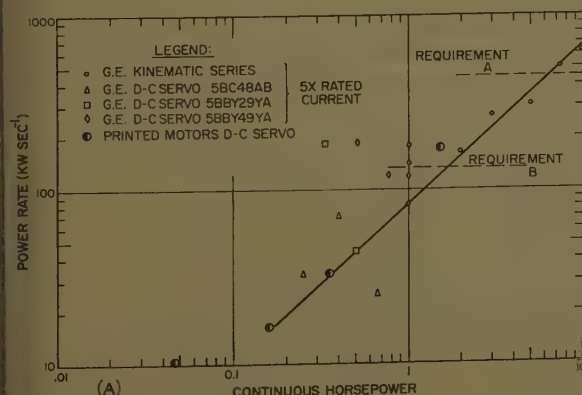


Fig. 2. Power rate versus (A) continuous power rating and (B) weight of typical servo-motors

Table II. Comparison of Actuators Meeting Requirements A

Motor	Power Rate, Kw/Sec <sup>-1</sup>	Weight, Pounds	I <sup>2</sup> R Loss, Kw
7 1/2-hp.....	500.....	240.....	7.30
10-hp.....	622.....	325.....	14.60

7 1/2- and 10-hp (horsepower) General Electric kinematic motors which are compared in Table II. For missile applications, actuator weight must be minimized, and thus, it is obvious that the 7 1/2-hp General Electric kinematic motor is the lightest one that will meet the dynamic performance specification.

The conclusion, drawn from this study, is that the motor weight is too great in relation to hydraulic apparatus. The combined weight of the three motors for the required three actuators, in this example, is 720 pounds exclusive of control circuitry and power source. A complete hydraulic system needed to perform this same function weighs approximately 50 pounds including control circuitry and power source. A lesser, though still significant, contrast in weight between electric and hydraulic systems exists for requirement B. This is because the 1/3-hp General Electric servomotor is the lightest motor that can meet the less stringent requirements and this device alone weighs 25 pounds. In view of these findings it is important to determine whether electric devices of suitable weight are, or are not, intrinsically capable of achieving the necessary power rates for missile applications.

This is discussed in the next section for a particular type of motor.

Theoretical Limit on Power Rate of Moving Conductor Actuators

Before rejecting electric motors for missile actuator applications it is desirable to explore the theoretical limit on the power rate of such devices as a function of size. At present it is impossible to discuss such performance limits in general terms; however, for particular types of electric motors limits on power rate can be determined. This section will consider motors with moving conductors of the d-c type. For d-c motors, the two following questions must be answered:

- 1. How should the motor be proportioned in order to produce the maximum power rate for given size?
- 2. What is the maximum possible power rate for given size, or conversely, what is the least size the motor may have for a given power rate?

The shell-type motor, shown in Fig. 3, consists of an armature in the form of a hollow shell which revolves around a stationary ferromagnetic core. Field flux through the poles is provided by an external magnetic circuit which is not shown. It is assumed that the magnetic circuit is capable of producing a field flux density,  $B_f$ , under the poles. The air gap, which is larger than in conventional motors, is assumed to be sufficiently small so that a substantially uniform flux density exists under the poles. The armature shell is assumed to have a mean resistivity,  $\rho_r$ , and a mean density,  $\rho_m$ . These values account for conductor insulation and structural materials used to bond the conductors into the shell form.

Utilizing these assumptions, the analysis in the Appendix shows that the power rate of the motor is proportional to the armature  $I^2R$  loss,  $P_R$ . Specifically

P\_M = (alpha\_p^2 / (alpha\_R alpha\_J)) \* (B\_f^2 / (rho\_M rho\_R)) \* P\_R (10)

where  $\alpha_p$  is the fraction of the shell under the field poles,  $\alpha_R$  the ratio armature total resistance to the shell resistance, and  $\alpha_J$  the ratio of armature total inertia to the shell inertia. Since these quantities cannot be varied appreciably from one motor design to another, it is concluded that the motor power rate is substantially independent of the motor dimensions for a given power loss in the armature.

In answer to the first question using equation 10, the motor proportions have very little influence on the power rate that can be had from a motor of a given size. This result contrasts sharply with the accepted practice in the design of servomotors having a small armature diameter-to-length ratio. The result of equation 10 shows that the important parameter defining the capability of a servomotor, namely, its power rate, is independent of the armature diameter to length ratio and depends only upon the power dissipated in  $I^2R$  losses in the armature.

Fig. 4 shows power rate as a function of armature  $I^2R$  loss for a number of commercially available d-c motors. The line drawn through the scattered points illustrates that the power rate tends to increase as the 1.105 power of the  $I^2R$  loss rather than as the first power as predicted by equation 10. However, in view of the variety of motors plotted in Fig. 4, the agreement is better than one could reasonably expect. It is also interesting to note that 46 seconds<sup>-1</sup> is the ratio,  $P_M/P_R$ , of the power rate to  $I^2R$  loss for a motor of a power rate of 137 kw/sec. This number is a measure of configuration efficiency

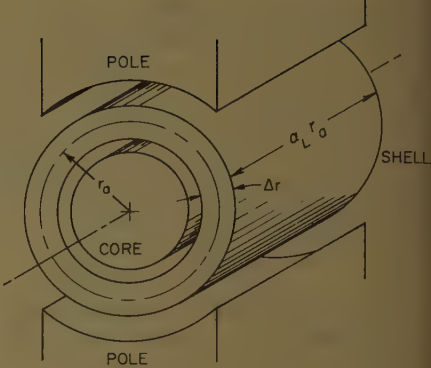


Fig. 3. Shell-type d-c motor

for a motor in producing power rate. For a shell-type motor, equation 10 shows that a ratio,  $P_M/P_R$ , of 724 seconds<sup>-1</sup> should be obtainable. This figure assumes that  $\alpha_p=2/3$ ,  $\alpha_R=2$ ,  $\alpha_J=2$ ,  $B_f=1.0$  weber meter<sup>-2</sup>,  $\rho_M=2.71 \times 10^8$  kg-meter<sup>-3</sup>, and  $\rho_R=5.66 \times 10^{-8}$  ohm-meter. The density corresponds to aluminum on the assumption that aluminum conductors and bonding materials of approximately equal density are used. The resistivity is twice that of aluminum on the assumption that only approximately one half of the shell cross section is conductor. An improvement in power rate capability for a given  $I^2R$  loss of more than 15 through configuration improvement certainly appears worth striving for.

The result of equation 10 indicates that there is no limit on the size of a motor for a given power rate, providing ways and means are found to dissipate the armature  $I^2R$  losses without excessive temperature rise. However, armature reaction must be considered. In the upper section of Fig. 5, the development of the air gap in the vicinity of one of the poles is shown. The lower portion shows the armature reaction mmf (magnetomotive force) as a function of the angle from the pole center. The peak armature reaction mmf occurs at the edge of the pole. As the dissipated power in the armature is increased, the armature reaction mmf at the pole edge will increase until it is excessively large in relation to the field mmf at this point. Although the armature reaction could be cancelled by pole face windings, this is not feasible in view of the additional  $I^2R$  losses.

In the absence of pole face windings the armature reaction mmf must be limited in relation to the field in mmf so that excessive distortion of the resultant flux pattern under the pole is avoided. The ratio of the maximum armature reaction mmf,  $F_{a(max)}$ , in Fig. 5, to the field mmf under the pole will be designated as



This ratio is considered to be a design specification and will have values ranging from 1/2 to 1 for motors with shell-type armatures under consideration. In the analysis in the Appendix it is shown that the constraint on the armature reaction mmf is equivalent to a constraint on the armature shell thickness,  $\Delta r$ . This constraint is

$$2 \left( \frac{\pi \alpha_p^2 \mu_0^2}{2 \alpha_R \alpha_L \rho_R n^2 p \alpha_a^2 \alpha_g^3 B_f^3} \right) P_{R(\max)} \quad (11)$$

where,  $\mu_0$  is the permeability of free space;  $\alpha_L$  is the ratio of air gap length to armature shell thickness;  $n_p$  is the number of poles;  $\alpha_L$  is the ratio of the armature shell length to its radius. All of the other symbols have been previously defined. Thus, by relation 11, a lower limit is placed on the thickness of the armature shell, for a specified armature  $I^2R$  loss,  $P_{R(\max)}$ , and a specified limitation,  $\alpha_r$ , is imposed on armature reaction.

Another constraint on the dimensions of the armature shell is imposed by cooling considerations. The armature shell must have a sufficiently large surface to transfer the  $I^2R$  losses to the cooling medium within the permissible temperature rise,  $\Delta\theta_{(\max)}$ . It is assumed that cooling occurs on the external surface of the armature shell. The maximum permitted  $I^2R$  loss must be equal to or greater than the average  $I^2R$  loss that actually occurs. Thus, the following equation can be written

$$U \alpha_s 2 \pi r_a^2 \Delta\theta_{(\max)} \geq \alpha_d P_{R(\max)} \quad (12)$$

where  $U$  is the over-all heat transfer coefficient;  $\alpha_s$  is the ratio of the total external armature surface to the shell surface;  $\Delta\theta_{(\max)}$  is the maximum permissible temperature rise of the armature shell relative to the cooling medium;  $r_a$  is the mean radius of the armature shell; and  $\alpha_d$  is the fractional duty cycle. This relation can be written as a constraint on armature radius as follows:

$$r_a \geq \sqrt{\frac{\alpha_d P_{R(\max)}}{2 \pi U \alpha_s \Delta\theta_{(\max)}}} \quad (13)$$

Equations 11 and 13 establish the dimensions

of the armature shell except for the ratio  $\alpha_L$  of the length to the radius. There may be some advantage, in equation 11, in making this ratio large in order to reduce the armature shell thickness and thus reduce the length of the air gap. A small air gap involves a small amount of permanent magnet material for the field (or field ampere-turns if electromagnets are used to establish the field flux). However, increasing the length of the armature, in relation to its radius, requires a longer field structure, and offsets, at least partially, the weight saving obtained through the narrower air gap. An exact evaluation of the optimum ratio of length to radius for the armature requires an analysis of the field design problem, which cannot be treated in this paper due to space restrictions. It will suffice to say that commonly used length-to-radius ratios in the vicinity of 2 probably are not too far from optimum.

Thus, the procedure for designing electric motors for missile actuator use is now clear. First, the actuator requirements are analyzed to arrive at a requirement on the motor power rate (torque-squared-to-inertia ratio) as given by relation 6. Next, relation 10 is used to determine the value of armature  $I^2R$  loss,  $P_R$ , that is necessary to obtain the requisite power rate. Third, using a selected value of armature shell length to radius ratio,  $\alpha_L$ , relation 11 yields the lower limit that the armature shell thickness may have without excessive armature reaction. Normally this lower limit is used unless it is prohibited by fabrication difficulties since small air gaps are desirable to minimize the weight of the field structure. Finally, the radius of the armature is established by means of relation 13. Obviously the size of the motor will depend upon the heat transfer coefficient used in relation 11. For missile applications it is necessary to obtain rather large heat transfer coefficients and this may be done by circulating a mist of liquid particles through the air gap. After the basic motor dimensions have been established,

the field structure may be designed using standard design procedures. The next section presents a proposed motor design that was developed by this technique.

## Design of a Servomotor to Meet Requirement A

In this section, an example of a specially designed shell-type servomotor, which was designed to meet missile actuator requirement A, is given. The parameters of the previous section are assumed. From equation 10 it was shown that  $\dot{P}_M/P_R$  is equal to 724 sec<sup>-1</sup>. The motor power rate needed to meet requirement A is 452 kw/sec<sup>-1</sup>. If a factor of safety is incorporated and the required power rate is taken as 500 kw/sec<sup>-1</sup>, the resulting maximum required armature  $I^2R$  loss,  $P_{R(\max)}$ , is

$$P_{R(\max)} = \frac{500}{724} \approx 0.7 \text{ kw} \quad (10A)$$

The minimum allowable shell thickness consistent with the armature reaction constraint now can be found by relation 11. Here, in addition to those parameters previously given:

Ratio of shell length to radius  $\alpha_L = 2$   
Number of poles  $n_p = 2$   
Ratio of maximum armature reaction mmf under a pole to the field mmf  $\alpha_a = 0.5$   
Ratio of air gap length to shell thickness  $\alpha_g = 1.25$

Thus, relation 11 gives

$$\Delta r \geq \left[ \frac{\pi (0.667)^2 (4\pi \times 10^{-7})^2}{(2)(2)(2)(5.66 \times 10^{-8})(2)^2 \times (0.5)^2 (1.25)^2 (1)^2} \right]^{1/2} 700 \quad (11A)$$

and

$$\Delta r_{(\min)} \approx 2 \times 10^{-3} \text{ meters}$$

The mean shell radius,  $r_a$ , is determined by cooling considerations. The following parameters are applicable:

Over-all heat transfer coefficient  $U = 340$  watts/m<sup>2</sup>C (spray mist cooling assumed)  
Duty cycle fraction  $\alpha_d = 0.25$

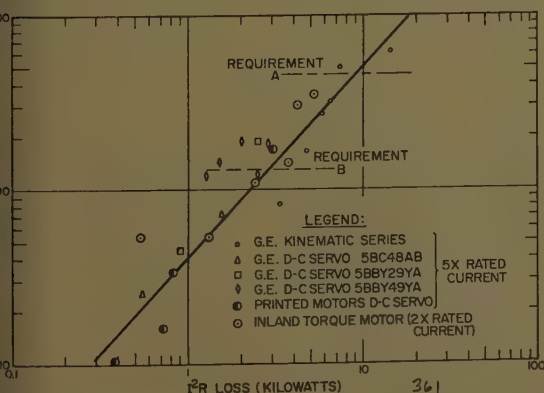


Fig. 4 (left). Power rate versus armature  $I^2R$  loss of typical servomotors

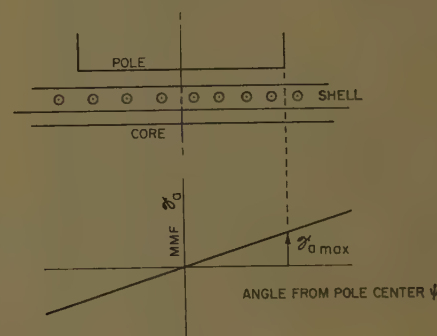


Fig. 5 (right). Armature reaction mmf

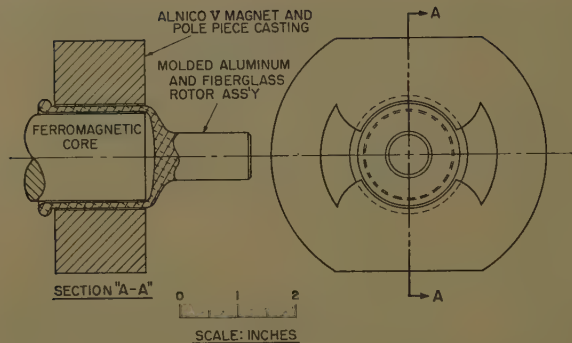


Fig. 6. Sketch of servomotor to meet requirement A

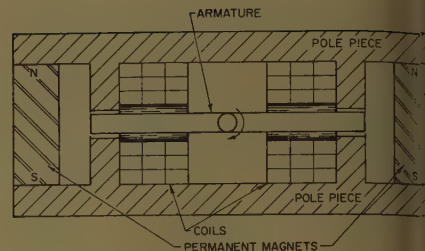


Fig. 7. Reluctance-type motor

Ratio of total armature surface to shell surface  $\alpha_s = 2$   
Maximum temperature rise  $\Delta\theta_{\max} = 50$  degrees centigrade

It should be noted that a 25% duty cycle,  $\alpha_d$  is used. This value greatly exceeds the 4% that was used in the comparison of conventional motors. With the indicated parameters, relation 13 gives

$$r_a \geq \sqrt{\frac{(0.25)(700)}{2\pi(340)(2)(2)(50)}} \quad (13A)$$

and

$$r_{a(\min)} \cong 0.02 \text{ meter}$$

Thus, the general motor configuration has been determined since shell length, air gap, and pole face width are all related to either  $\Delta r_a$  or  $r_a$  by the various design ratios.

A sketch of the rotor and pole piece of the motor is shown in Fig. 6. It can be shown, by conventional analysis, that the permanent magnet configuration will provide the required 1 weber-per-square-meter flux density (10,000 gauss). With 200 armature conductors the brush current will be approximately 31.2 amperes at maximum torque.

The motor, illustrated in Fig. 6, has the following characteristics:

Shell moment of inertia  $J_m = 2.18 \times 10^{-5}$  kg-m<sup>2</sup>  
Maximum motor torque  $T_m = 3.32$  newton-meters

The resulting power rate is the required value of 505 kw/sec<sup>-1</sup>. Therefore, the motor design is consistent with the load requirements. The ferromagnetic material of the design weighs only 7 pounds and the estimated weight of the motor, exclusive of cooling equipment, is less than 11 pounds. In a 3-actuator system, the weight of the cooling system and control circuitry should be less than the weight of the motors alone. Thus, an over-all system weight of 50 to 60 pounds should be achievable. This is exclusive of the increment in power supply weight needed to supply the required electric power. Since the hydraulic system

needed to meet these requirements weighs about 50 pounds it is evident that electric actuators may be feasible in this application.

It should be noted that if a 2.5% duty cycle were assumed, no spray mist cooling would be required for the proposed design. Also, requirement A is representative of a rather severe set of actuator requirements. With a reduction in duty cycle or performance requirements, the weight of an electric actuator system can be further reduced and its competitive position in relation to hydraulic systems can thereby be improved.

### Prospects for Reluctance-Type Motors

In the preceding section a specially designed, moving conductor, d-c motor was shown to have reasonably good prospects for meeting the power rate and weight requirements that are typical of missile actuator applications. The question arises whether this is the only type of motor that shows such promise. A partial answer to this question is supplied by briefly examining reluctance-type motors in which torque or force is developed on the moving element because of a changing reluctance of an air gap. A relay is an elementary example. Fig. 7 shows a more sophisticated type of reluctance motor. An analysis of the theoretical limit on the power rate of this type of motor, similar to that conducted for moving conductor devices, shows that the ratio of power rate to  $I^2R$  loss ranges from 10,000 to 50,000 inverse seconds. These figures are substantially larger than for moving conductor devices which range from 100 to 1,000 inverse seconds.

In order to lend credence to these figures, data are cited in Table III for a commercially available torque motor for actuating hydraulic control valves that utilizes the configuration of Fig. 7.

Unfortunately, the favorable power rate to  $I^2R$  ratios for reluctance-type motors are obtained only for small air gaps, and therefore, for limited linear or

angular motion. To obtain continuous rotation it is necessary to employ a plurality of motors of this type in conjunction with a mechanism for converting oscillatory motion into continuous rotation. This mechanism would be analogous to connecting rod and crank mechanism used in reciprocating engines. If the inertia introduced by the coupling mechanism is negligible in relation to the basic motor inertias, it can be shown that for two motors moving sinusoidally 90 degrees out of phase, the resultant torque squared-to-inertia ratio or power rate at the output shaft will be equal to that characterizing the individual motors. Thus, assuming that suitable coupling mechanisms or linkage can be found, the reluctance-type motor should be able to provide the basis for an actuator of power rate capability that would be quite satisfactory for missile applications.

Another disadvantage of the reluctance-type motor is the necessity of supplying the individual motor elements with alternating currents of adjustable frequency. In effect, a plurality of such motor elements driving a common shaft in the manner described constitutes a synchronous motor. For a typical missile application, the frequencies of the supplied currents will range from zero to several hundred cycles per second. The complication of a variable-frequency source that is capable of supplying controlled currents against back electromotive forces of amplitude that are proportional to frequency, constitutes a serious handicap for this type of motor.

Table III. Data on Torque Motor

Manufacturer	Midwestern Geophysical Laboratory
Model no.	9
Moment of inertia $J_M$	$64 \times 10^{-6}$ inch-pound-sec <sup>2</sup>
Peak torque $T_M$	8.33 inch-pounds
Power rate $\dot{P}_M$	132 kw sec <sup>-1</sup>
Peak current $i_c$	40 milliamperes
Coil resistance $R_c$	3,400 ohms
Coil inductance $L_c$	7 henrys
$I^2R$ loss $P_R$	5.44 watts
Ratio $\frac{\dot{P}_M}{P_R}$	$24.3 \times 10^3$ sec <sup>-1</sup>
Weight	1.2 pounds



Considering both the complications of the mechanism required to convert oscillatory to rotary motion and the required variable-frequency current source, the prospects for reluctance-type motors for missile actuator applications appear to be questionable, in spite of the large power rates. However, the harmonic-drive concept<sup>6</sup> may provide an answer to the problem of converting oscillatory to rotary motion, and therefore, in view of their large theoretical power rates, research will continue on actuators using reluctance-type motors.

## Conclusions

This paper has shown that the torque-squared-to-inertia ratio or power rate of a motor is the most important characteristic for determining its applicability as a missile control surface actuator. Examination of typical missile control surface requirements shows that the required power rates range from 100 to 500 kw per second.

An examination of conventional motors of both the ordinary and printed circuit variety has shown that they are incapable of meeting the power rates required in missile applications within sizes and weights that are useful. However, an analysis of the theoretical limit on the power rates that can be achieved by moving conductor d-c motors shows that the power rate tends to be proportional to the  $T^2/R$  losses in the armature and independent of the motor size. The dimensions of the armature tend to be fixed by armature reaction and cooling considerations. Using these results, a specially designed d-c motor of reasonable weight and size was shown to be capable of meeting the imposed power rate requirements of a rather difficult missile application. Finally, another type of electric actuator in the form of a reluctance-type motor has been examined. Although this type of motor is capable of rather extreme power rates, it suffers from two serious disadvantages: the need for special mechanisms for converting oscillatory motion into rotation, and the need for a power supply that is adjustable in both frequency and current magnitude.

This paper has shown that it is theoretically feasible to develop electric actuators for missile applications insofar as dynamic performance and size and weight requirements are concerned. The overall systems problem including electric power sources and controls has not been

discussed, but this must be considered in each application before deciding whether or not to use electric actuators.

## Appendix

The derivation of equation 8 is as follows. The basic expression for the torque developed by the armature is

$$T_M = \int dv [\mathbf{r} \times (\mathbf{j} \times \mathbf{B}_f)] \quad (14)$$

where  $\mathbf{j}$  is the current density in the armature shell and  $\mathbf{r}$  is the radius vector to the volume element,  $dv$ , from a point on the axis of rotation. The other symbols have been previously defined; all symbols are defined in the nomenclature. The integration is to be carried out over all of the armature volume. Using the previously defined nomenclature the magnitude of the torque turns out to be

$$T_M = 2\pi\alpha_p\alpha_L r_a^3 \Delta r j B_f \quad (15)$$

The moment of inertia of the armature, assuming that there is a thin shell, is given by

$$J_M = 2\pi\alpha_J\alpha_L\rho_m r_a^4 \Delta r \quad (16)$$

so that the torque-squared-to-inertia ratio is

$$\frac{T_M^2}{J_M} = \frac{2\pi\alpha_p^2\alpha_L^2 r_a^2 \Delta r j^2 B_f^2}{\alpha_J\rho_m} \quad (17)$$

The power dissipated in the armature is given by the volume integral

$$P_R = \int dv (j^2 \rho_R) \quad (18)$$

For the shell-type armature under consideration the integration yields

$$P_R = 2\pi\alpha_R\alpha_L\rho_R r_a^2 \Delta r j^2 \quad (19)$$

Using this equation to eliminate  $j$  from equation 17 yields equation 10 since the torque-squared-to-inertia ratio is defined as the power rate  $\dot{P}_M$ .

The constraint on the armature shell thickness  $\Delta r$  given by relation 11 is derived as follows. Referring to Fig. 5 the basic armature reaction constraint is

$$\mathcal{F}_{a(\max)} \geq \alpha_a \mathcal{F}_f \quad (20)$$

which states that the maximum armature reaction mmf,  $\mathcal{F}_{a(\max)}$ , under a pole shall be equal to or less than the field mmf,  $\mathcal{F}_f$ , multiplied by a specified factor  $\alpha_a$ . In terms of the field flux density,  $B_f$ , and the armature shell thickness this expression becomes

$$\mathcal{F}_{a(\max)} \leq \alpha_a \frac{B_f}{\mu_0} \alpha_g \Delta r \quad (21)$$

since the air gap is assumed to be proportional to  $\Delta r$ . From Fig. 5 it is evident that the armature reaction mmf at the edge of the pole is given by

$$\mathcal{F}_{a(\max)} = \frac{\pi\alpha_p r_a \Delta r j_{(\max)}}{n_p} \quad (22)$$

Squaring both sides and using equation 19 to eliminate  $j$  yields

$$\mathcal{F}_{a(\max)}^2 = \frac{\pi\alpha_p^2 \Delta r P_{R(\max)}}{2n_p^2 \alpha_R \alpha_L \rho_R} \quad (23)$$

Squaring both sides from equation 21 and using the value of  $\mathcal{F}_{a(\max)}$  given by equation 23 yields equation 11 upon rearrangement.

## Nomenclature

- $a_L$  = peak load acceleration
- $R_G$  = ratio of gearing
- $v_L$  = peak load velocity
- $v_M$  = peak motor velocity
- $J_L$  = moment of inertia of load
- $J_M$  = moment of inertia of motor
- $M_M$  = mass of motor
- $\dot{P}_M$  = peak power rate of motor
- $T_L$  = peak torque of load
- $T_M$  = peak torque of motor
- $P_R = I^2 R$  loss in armature
- $r_a$  = mean radius of armature shell
- $\rho_R$  = mean resistivity of armature shell
- $\rho_m$  = mean density of armature shell
- $B_f$  = mean field flux density under pole
- $\alpha_p$  = fraction of shell under poles
- $\alpha_J$  = ratio of total inertia to that of armature shell
- $\Delta r$  = thickness of armature shell
- $\alpha_R$  = ratio of armature total resistance to that of shell
- $n_p$  = number of poles
- $\alpha_R$  = ratio of shell length to radius
- $\mu_0$  = permeability of free space
- $\alpha_a$  = ratio of maximum allowable armature reaction mmf under pole to field mmf under pole
- $\alpha_g$  = ratio of air gap length to armature shell thickness
- $\Delta\theta$  = temperature rise of armature shell
- $\alpha_s$  = ratio of total external armature surface to shell surface
- $\alpha_d$  = fractional duty cycle
- $U$  = over-all heat transfer coefficient
- $\mathcal{F}_a$  = armature reaction mmf
- $\mathbf{j}$  = current density in armature shell
- $\mathcal{F}_f$  = field mmf under pole

## References

1. FLUID POWER CONTROL (book), John F. Blackburn, Gerhard Reethof, J. Lowen Shearer, editors. Technology Press, Cambridge, Mass.; also John Wiley & Sons, Inc., New York, N. Y., 1960.
2. DEVELOPMENT OF A GAS-OPERATED REACTION-JET SERVOMOTOR, R. S. Scher. *Proceedings, British Institution of Mechanical Engineers, "Symposium on Automatic Control,"* London, England, 1960.
3. APPLICATION OF SILICON-CONTROLLED RECTIFIERS IN A TRANSISTORIZED HIGH-RESPONSE D-C SERVO SYSTEM, Clarence Cantor. *AIEE Transactions*, pt. II (*Applications and Industry*), vol. 80, Mar. 1961, pp. 7-12.
4. WHAT SIZE MOTOR FOR PROPER OPERATION OF SERVOMECHANISM, G. C. Newton. *Machine Design*, Cleveland, Ohio, vol. 22, Nov. 1950, pp. 125-30.
5. A COMPARISON OF TWO BASIC SERVOMECHANISM TYPES, Herbert Harris. *AIEE Transactions*, vol. 66, 1947, pp. 83-93.
6. BREAKTHROUGH IN MECHANICAL DRIVE DESIGN: THE HARMONIC DRIVE, C. W. Musser. *Machine Design*, vol. 32, Apr. 14, 1960, pp. 160-7.

(See page 312 for discussion)



## Discussion

N. F. Tsang (University of Arkansas, Fayetteville, Ark.): The authors have given an interesting approach to the solution of high performance actuators. In the specially designed shell-type servomotor it may be noted that a very important design parameter is  $B_f$ , the air gap flux density. The allowable amount of total ampere conductors of the armature is proportional to the air gap magnetomotive force, which is in turn proportional to  $B_f$ . Hence the torque will be proportional to  $B_f^2$ , and  $P_m$  to  $B_f^4$ . A 10% reduction in  $B_f$  would mean 34% in  $P_m$ , the performance factor.

It would be enlightening if the authors would give a more detailed analysis of the magnetic circuit of the specially designed motor, stressing the following points:

1. Shape and magnetization of the permanent magnet.
2. Leakage fluxes.
3. Flux reduction due to field distortion.

It seems somewhat doubtful that the assumed  $B_f$  of 1 weber-meter<sup>2</sup> could be attained with Alnico V in the design shown. Even though it cannot be realized, some other material, such as Ticonal or Alcomax, may be considered, or the magnetic circuit reshaped, or both.

G. C. Newton, Jr., R. W. Rasche: The shape and size of the permanent magnet

shown is intended to represent a minimum size permanent magnet for the device proposed. We recognize that some pole shaping or other means of field compensation or concentration might be required. It is felt, however, that such changes would not result in any appreciable size or weight change. Since the application considered in the paper was an extremely stringent one, the conclusions of the paper would not be affected by minor dimensional changes in the device proposed, or for that matter, by an appreciable degradation of the performance characteristics of the actuator.

The magnetic circuit analysis is based on the assumption that the mmf producing the air gap flux is developed primarily in the side legs of the magnet. The additional magnetic material will only increase air gap flux. The cross-sectional area of each leg is  $0.6 A_g$  where  $A_g$  is the cross-sectional area of the air gap under the pole face. A leakage factor of 1.2 is assumed. Fringing is neglected.

Using

$$B_m = \frac{f B_f A_g}{A_m} \quad (24)$$

where  $B_m$  is magnet flux density in side legs,  $f$  the leakage factor,  $A_g$  the cross-sectional area of air gap under pole face, and  $A_m$  the cross-section area of magnet side legs, and realizing that, by design,

$$A_m = 1.2 A_g \quad (25)$$

The required flux density is

$$B_m = \frac{(1.2)(1.0)(A_g)}{1.2(A_g)} = 1.0 \text{ weber-meter}^{-2} \quad (26)$$

Reference to  $B$ - $H$  curves for Alnico V shows 1.0 weber-meter<sup>-2</sup> corresponds to  $H_m$  equal to 35,000 ampere turns per meter<sup>-1</sup>.

The total gap width is

$$2\alpha_g \Delta r = 0.005 \text{ m} \quad (27)$$

The required length of magnet is

$$L_m = \frac{B_f 2\alpha_g \Delta r \times 10^7}{4\pi H_m} \quad (28)$$

or

$$L_m = \frac{(1.0)(.005) \times 10^7}{4\pi(35,000)} = 0.11 \text{ m} \quad (29)$$

The effective length of each side leg is at least 0.11 m depending on how much of the pole pieces are considered. Thus it is felt that the resultant gap flux density will sufficiently exceed 1.0 weber-meter<sup>-2</sup> to allow for distortion effects. Additional material in the side legs could improve this figure by 25%.

It should be noted that equation 11 is the result of an armature reaction constraint. This constraint can be given any value, but it is our opinion that a value of  $\alpha_g$  of 0.5 will result in sufficiently low distortion to assure satisfactory operation.

The suggestions of Mr. Tsang as to the use of Ticonal or Alcomax are valid and are much appreciated.



## Power Apparatus and Systems—October 1961

61-204	Surge Comparison Testing of D-C Armature Windings.....	Scheda . . .	537
61-143	Collector Ring and Brush Wear.....	Herder, Kerber . . .	543
61-762	Basic Constants of General Induction Machine.....	Graybeal . . .	548
61-807	Repeater Receiver for the NEAR System.....	Cleary . . .	556
61-633	32-Step Voltage-Regulator Performance.....	Gangel, Green . . .	559
61-760	Self-Compensated High-Current Homopolar Generator....	Das Gupta . . .	567
61-720	Photoelectric Impulse Generation for Demand Metering....	Whipple . . .	573
61-173	High-Frequency Iron Losses in Fractional-Hp Motors.....	Shaneman . . .	579
61-740	Field Patterns of Bundle Conductors.....	Timascheff . . .	590
61-761	Saturation Harmonics of Polyphase Induction Machines.....	Lee . . .	597
61-489	Mechanism of Breakdown of Laboratory Gaps.....	Wagner, Hileman . . .	604
61-488	The Lightning Stroke—II.....	Wagner, Hileman . . .	622
61-737	Protection of Lead-Sheathed Power Cables.....	Trouard, Maier . . .	642
61-238	Trends in H-V Cable Techniques in Great Britain.....	Barnes, Sutton . . .	647
60-832	Determination of Economical Distribution Substation Size.....	Smith . . .	663
60-833	Economics of Primary Distribution Voltages of 4.16–34.5 Kv.....	Smith . . .	670
61-757	Simplified Methods of Calculating Insulation Life.....	Whitman . . .	683
61-772	The Cowans Ford Project.....	Wray . . .	685

### Conference Papers Open for Discussion

Conference papers listed below have been accepted for AIEE Transactions and are now open for written discussion until December 27. Duplicate double-spaced typewritten copies for each discussion should be sent to Edward C. Day, Assistant Secretary for Technical Papers, American Institute of Electrical Engineers, 345 East 47th Street, New York 17, N. Y., on or before December 27.

Preprints may be purchased at 50¢ each to members, \$1.00 each to non-members, if order is accompanied by remittance or coupons. Please order by number and send remittance to:

AIEE Order Department  
345 East 47th Street  
New York 17, N. Y.

60-868	Stabilization of Fluid Servomechanisms by the Stability Factor Method.....	Fonda
61-85	Servomotor Characteristics by Impulse Testing.....	Weed
61-223	High-Accuracy Digital-Analog Solid-State Speed Controller....	Thompson
61-645	Speed Regulation by Digital Methods.....	Potts
61-887	Application and Performance of Insulations for 500 C Hypersonic Aircraft Generators.....	Balke, Merrifield, Dimond
61-889	Optimizing Simple Circuitry for Reliability and Performance by Failure Mode.....	Hanne
61-891	Static Inverter with Neutralization of Harmonics.....	Kernick, Roof, Heinrich
61-906	Problems in the Selection and Testing of Nickel Cadmium Batteries for Satellites.....	Albrecht
61-911	Solar-Cell Performance with Concentrated Sunlight.....	Tallent, Oman
61-913	Voltage Modulation on Aircraft Electric Power Systems as a Function of the Flexibility of the Alternator Drive Shaft.....	Flugstad



# AIEE PUBLICATIONS

Member Prices	Nonmember Prices	
	Basic Prices*†	Extra Postal for Foreign Subscription

## Electrical Engineering

Official monthly publication containing articles of broad interest, technical papers, digests, and news sections: Institute Activities, Current Interest, New Products, Industrial Notes, and Trade Literature. Automatically sent to all members and enrolled students in consideration of payment of dues. (Members may not reduce the amount of their dues payment by reason of nonsubscription.) Additional subscriptions are available at the nonmember rates.

annually	\$12*	\$1.00
Single copies	\$1.50*	

## Bimonthly Publications

Containing all officially approved technical papers collated with discussion (if any) in three broad fields of subject matter as follows:

	annually	annually	
Communication and Electronics	\$5.00	\$8.00*	\$0.75
Applications and Industry	\$5.00	\$8.00*	\$0.75
Power Apparatus and Systems	\$5.00	\$8.00*	\$0.75

Each member may subscribe to any one, two, or all three bimonthly publications at the rate of \$5.00 each per year. A second subscription to any or all of the bimonthly publications may be obtained at the nonmember rate of \$8.00 each per year.

Single copies may be obtained when available.	\$1.50 each	\$1.50* each
---	----------------	-----------------

## AIEE Transactions

An annual volume in three parts containing all officially approved technical papers with discussions corresponding to six issues of the bimonthly publication of the same name bound in cloth with a stiff cover.

	annually	annually	
Part I Communication and Electronics	\$4.00	\$8.00*	\$0.75
Part II Applications and Industry	\$4.00	\$8.00*	\$0.75
Part III Power Apparatus and Systems	\$4.00	\$8.00*	\$0.75

Annual subscription to all three parts (beginning with vol. 77 for 1958).

\$10.00	\$20.00*	\$2.25
	\$15.00*	\$1.50

Annual subscription to any two parts.

## AIEE Standards

Listing of Standards, test codes, and reports with prices furnished on request.

## Special Publications

Committee reports on special subjects, bibliographies, surveys, and papers and discussions of some specialized technical conferences, as announced in ELECTRICAL ENGINEERING.

\*Discount 25% of basic nonmember prices to college and public libraries. Publishers and subscription agencies 15% of basic nonmember prices. For available discounts on Standards and special publications, obtain price lists from Order Department at Headquarters.

†Foreign prices payable New York exchange

Send all orders to:

Order Department  
American Institute of Electrical Engineers  
345 East 47th Street, New York 17, N. Y.